

Redes neuronales y aprendizaje por transferencia



UNIVERSIDAD COMPLUTENSE DE MADRID

FACULTAD DE CIENCIAS MATEMÁTICAS

Trabajo de Fin de Grado

Curso: 2022 - 2023

Grado en Matemáticas

Autor

FRANCISCO JAVIER ABOLLADO MÚGICA

Tutor

ANTONIO LÓPEZ MONTES

Resumen

Este trabajo se estructura en dos secciones principales que abordan los fundamentos y aplicaciones contemporáneas de las redes neuronales. La primera parte, **Introducción a las redes neuronales**, proporciona un panorama exhaustivo sobre el desarrollo inicial y las estructuras de las redes neuronales, incluyendo el perceptrón, el Perceptrón Multicapa (MLP), las Redes Neuronales Convolucionales (CNN) y las Redes Neuronales Recurrentes (RNN). Se detallan además los componentes clave como las funciones de activación, las funciones de pérdida y los optimizadores. Esta sección culmina con ejemplos prácticos de implementación de redes neuronales en Python, consolidando así la base teórica con la aplicación práctica.

La segunda parte del trabajo se enfoca en **Aprendizaje por transferencia**, una metodología avanzada y eficiente para la aplicación de redes neuronales en problemas de la vida real. Aquí, se discuten los beneficios significativos del aprendizaje por transferencia, explicando cómo se puede mejorar el rendimiento de un modelo mediante la reutilización de conocimientos preaprendidos de conjuntos de datos extensos como ImageNet. Se presentan ejemplos prácticos en python de cómo esta técnica se aplica en modelos avanzados, como Distilbert, para ilustrar su efectividad y versatilidad en diversas aplicaciones. Este segmento destaca la transición de la teoría a la práctica, enfatizando la relevancia y la potencialidad del aprendizaje por transferencia en el ámbito de la inteligencia artificial.

Palabras clave

Redes neuronales, aprendizaje automático, clasificación, aprendizaje por transferencia, ImageNet, DistilBERT, python

Abstract

This pdf is structured into two main sections that explore the foundations and contemporary applications of neural networks. The first part, **Introduction to Neural Networks**, provides a comprehensive overview of the initial development and structures of neural networks, including the perceptron, Multi-Layer Perceptron (MLP), Convolutional Neural Networks (CNNs), and Recurrent Neural Networks (RNNs). It also details key components such as activation functions, loss functions, and optimizers. This section concludes with practical examples of implementing neural networks in Python, thereby solidifying the theoretical foundation with practical application.

The second part of the thesis focuses on **Transfer Learning**, an advanced and efficient methodology for applying neural networks to real-world problems. Here, the significant benefits of transfer learning are discussed, explaining how model performance can be enhanced by reusing knowledge pre-learned from extensive datasets like ImageNet. Practical examples in python of how this technique is applied in advanced models, such as Distilbert, are presented to illustrate its effectiveness and versatility across different applications. This segment highlights the transition from theory to practice, emphasizing the relevance and potential of transfer learning in the field of artificial intelligence.

Keywords

Neural networks, machine learning, classification, transfer learning, ImageNet, DistilBERT, python

Índice

Resumen	2
Abstract	3
Índice de figuras	7
I PARTE 1: Introducción a las redes neuronales.	8
1. Red neuronal básica. El perceptrón	9
1.1. Neurona artificial	9
1.2. Arquitecturas previas al perceptrón	10
1.3. Historia del perceptrón	11
1.4. Concepto de entrenamiento de un red neuronal	11
2. Tipos de arquitecturas en las redes.	13
2.1. Perceptrón simple	13
2.2. Perceptrón multicapa	14
2.3. Red convolucional	15
2.3.1. Capa convolucional	15
2.3.2. Capas de Pooling	16
2.4. Red recurrente	17
2.4.1. Celdas de Memoria en RNN	17
3. Algunos elementos fundamentales.	19
3.1. Funciones de activación	19

3.1.1.	Unidad Lineal Rectificada (ReLU)	19
3.1.2.	SoftMax	19
3.1.3.	Sigmoide	20
3.1.4.	Tangente hiperbólica	20
3.2.	Función de Pérdida	21
3.2.1.	Ejemplos comunes	21
3.3.	Optimizadores	22
3.3.1.	Descenso del Gradiente Estocástico (SGD)	22
3.3.2.	Estimación adaptativa del momento (ADAM)	24
4.	Ejemplos de redes.	26
4.1.	Perceptrón simple - Separación lineal	26
4.2.	Redes Convolucionales - datos MNIST	28
4.3.	Redes Recurrentes para Predicción de Series Temporales	30
II	PARTE 2: Aprendizaje por transferencia.	34
5.	Introducción al apredizaje por transferencia	35
5.1.	Introducción	35
5.2.	Fundamentos	35
6.	Ejemplos de aprendizaje por transferencia	37
6.1.	Clasificación de Imágenes con ImageNet: Gatos vs Perros	37
6.1.1.	Configuración y Preparación de Datos	37
6.1.2.	Preparación del Modelo y Aprendizaje por Transferencia	38
6.1.3.	Entrenamiento del Modelo	39
6.1.4.	Evaluación y Conclusión	40
6.2.	Análisis de sentimientos con BERT	41
6.2.1.	Configuración y Preparación de Datos	41
6.2.2.	Preparación del Modelo y Aprendizaje por Transferencia	41
6.2.3.	Entrenamiento del Modelo BERT	42
6.2.4.	Evaluación y Conclusión	43

III Bibliografía y Anexos	45
Anexos	46
C. Backpropagation	46
C.1. Algoritmo	47
D. Aclaración de $L(\Theta)$	49
E. Clustering	49
E.1. Algoritmos de Clustering	50
E.2. K-means	50
E.3. DBSCAN	50
E.4. Evaluación de Clustering	51
Códigos	52
F. Perceptrón simple - Separación lineal	52
F.1. main.py	52
F.2. perceptron.py	53
G. Redes Convolucionales - datos MNIST	54
G.1. main.py	54
G.2. create_model.py	54
G.3. test_model.py	56
G.4. nn.py	56
G.5. files.py	58
G.6. ejemplo_layer1.py	58
H. Redes Recurrentes para Predicción de Series Temporales	59
H.1. main.py	59
I. Clasificación de Imágenes con ImageNet: Gatos vs Perros	61
I.1. main.py	61
I.2. test.py	63
J. Análisis de sentimientos con BERT	65
J.1. sentimental_analysis.ipynb	65
Bibliografía	71

Índice de figuras

1.1.1.Neurona Humana	9
1.1.2.Neurona Artificial	9
1.2.1.Clustering en 5 grupos	10
1.4.1.Entrenamiento de una red	11
2.1.1.Perceptrón de una capa	13
2.2.1.Perceptrón multicapa	14
2.3.1.Estructura de red convolucional	16
2.3.2.Pooling	16
2.4.1.Estructura de un red neuronal recurrente	17
2.4.2.Distintas estructuras de redes recurrentes	17
3.1.1.Relu	19
3.1.2.Softmax	20
3.1.3.Sigmoide	20
3.1.4.Tangente Hiperbólica	21
4.1.1.Evolución del perceptrón durante el entrenamiento	27
4.2.1.Filtros de la primera capa de la red	29
4.2.2.Ejemplo MNIST	30
4.2.3.Ejemplo MNIST tras el primer filtro	30
4.2.4.Ejemplo MNIST tras el primer filtro (V2)	30
6.1.1.Predicciones: Gato vs Perro	40
6.2.1.Matriz de confusión	44

Parte I

PARTE 1: Introducción a las redes neuronales.

Capítulo 1

Red neuronal básica. El perceptrón

1.1. Neurona artificial

Una red neuronal es un modelo matemático que emula el modo en que el cerebro humano procesa la información. La estructura de la red neuronal está formada por pequeños procesadores de información, llamados neuronas artificiales, basándose en el modelo de neuronas del cerebro humano desarrollado en sus inicios por Ramon y Cajal. De forma breve las similitudes entre la neurona humana y la artificial son las siguientes.

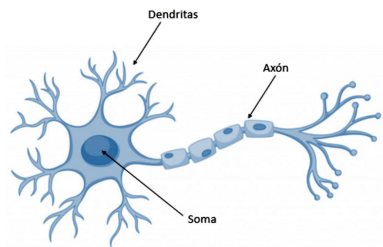


Figura 1.1.1: Neurona Humana

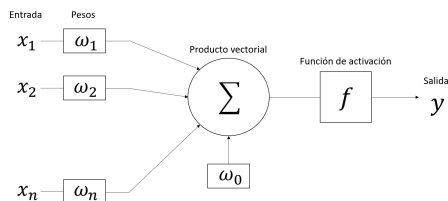


Figura 1.1.2: Neurona Artificial

Aspecto	Neurona Humana	Neurona Artificial
Canal de entrada de información	Las dendritas	Pesos de la neurona, denotados como $\vec{w} = (w_0, w_1, \dots, w_n) \in \mathbb{R}^{n+1}$. La información de entrada es un vector $\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$.

Órgano de cómputo	El soma	Producto vectorial entre \vec{w} y \vec{x} . Este resultado es procesado por una función de activación $f : \mathbb{R} \rightarrow \mathbb{R}$.
Canal de salida	El axón	El resultado del cómputo es guardado en la variable y .

Cuadro 1.1: Comparación entre neuronas humanas y artificiales

La estructura de la red neuronal se divide en tres capas principales: la capa de entrada (input layer), una o más capas intermedias (hidden layers) y la capa de salida (output layer). En esta estructura, las funciones de activación son una parte importante y se explicarán en [sección 3.1]. Su elección depende principalmente de dos factores: su derivabilidad y su coste computacional.

1.2. Arquitecturas previas al perceptrón

Previo al desarrollo de arquitecturas matemáticas más avanzadas como el perceptrón, encontramos el concepto de *clustering* [E], una técnica que parece compartir similitudes en términos de estructura y objetivo de resolución. El clustering es un enfoque de clasificación simple basado en la agrupación de datos en función de sus distancias, utilizando, por ejemplo, la distancia euclidiana.

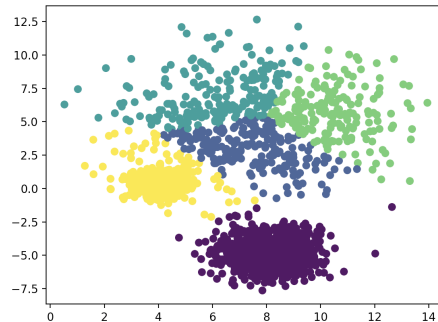


Figura 1.2.1: Clustering en 5 grupos

Clustering y la distancia euclidiana

El clustering se basa en la idea fundamental de *distancia*, un concepto matemático que tiene sus orígenes en la geometría euclidiana. Mediante la minimización de estas distancias, el clustering agrupa datos que están cercanos entre sí, proporcionando una forma básica de clasificación. Sin embargo, su utilidad está limitada a problemas lineales y de bajas dimensiones debido a su dependencia exclusiva del concepto de distancia.

Limitaciones del clustering y la necesidad de estructuras complejas

Aunque el clustering es eficaz para tareas simples de clasificación, su capacidad se ve restringida cuando se enfrenta a problemas más complejos y de alta dimensionalidad. La falta de una estructura subyacente y de mecanismos de adaptación impide que el

clustering aborde adecuadamente estos desafíos. Esto resalta la necesidad de modelos más sofisticados que puedan manejar la no linealidad y la mayor complejidad de los datos.

1.3. Historia del perceptrón

El perceptrón, una pieza clave en la historia de las redes neuronales, fue creado por Frank Rosenblatt en enero de 1957 en el Cornell Aeronautical Laboratory, ubicado en Buffalo, Nueva York. Rosenblatt y su equipo se propusieron crear una máquina capaz de aprender y reconocer patrones complejos, una iniciativa pionera que buscaba imitar funciones humanas avanzadas como la percepción visual, la formación de conceptos y la capacidad de generalizar a partir de experiencias previas.

El modelo básico del perceptrón es un clasificador lineal, lo que significa que solo puede trazar una línea de decisión para separar dos clases de datos. Funciona bien cuando los conjuntos de datos son linealmente separables, pero sus limitaciones se hacen evidentes con patrones más complejos que no se pueden separar simplemente con una línea.

Esta limitación fue destacada por Marvin Minsky y Seymour Papert en 1969 en [13], donde analizaron en profundidad las capacidades y las limitaciones del perceptrón, demostrando matemáticamente que no podía resolver problemas no lineales. Esto llevó a una reducción en el interés y la financiación de la investigación en redes neuronales hasta que en los 80 se introdujo las arquitecturas multicapa y el algoritmo de retropropagación. Estos avances revitalizaron el campo, permitiendo que las redes neuronales abordaran una gama más amplia y compleja de problemas de clasificación y regresión.

1.4. Concepto de entrenamiento de una red neuronal

Cuando se entrena una red neuronal, el objetivo principal es reducir al mínimo el error. Este error [sección 3.2] es la diferencia entre lo que la red predice y lo que realmente debería predecir según los datos de entrenamiento. Para entenderlo mejor, pensemos en la red neuronal como una herramienta que intenta aprender de ejemplos específicos.

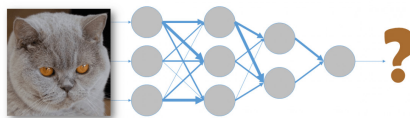


Figura 1.4.1: Entrenamiento de una red

Imagina que cada ejemplo de entrenamiento es un par de datos, representado como (\vec{x}, y) . Aquí, \vec{x} es un conjunto de características de entrada (como números o valores que describen algo específico), y y es el resultado correcto o esperado para esas características. Por ejemplo, en una tarea de clasificación, y podría ser la categoría a la que pertenece \vec{x} , y en una tarea de regresión, podría ser un valor numérico específico asociado a \vec{x} .

Cuando se introduce el conjunto de características \vec{x} en la red neuronal, la red da como resultado una predicción, que llamaremos z . Esta predicción es el intento de la red de adivinar el valor de y basándose en lo que ha aprendido hasta el momento.

El proceso de entrenamiento ajusta constantemente lo que la red ha aprendido con el objetivo de hacer que z (la predicción de la red) se acerque lo más posible a y (el resultado correcto). Esto se hace modificando los "pesos" de la red, que son parámetros que influyen en cómo la red procesa las entradas. Al ajustar estos pesos, la red se vuelve gradualmente mejor en hacer predicciones precisas.

Capítulo 2

Tipos de arquitecturas en las redes.

2.1. Perceptrón simple

El perceptrón simple, un modelo fundamental en la evolución de las redes neuronales artificiales, se distingue por su arquitectura minimalista capaz de efectuar tareas de clasificación en espacios linealmente separables.

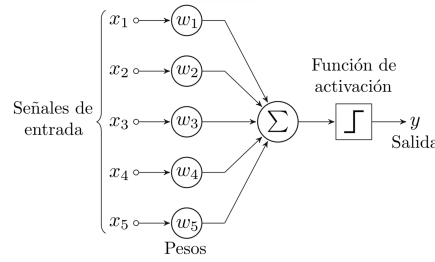


Figura 2.1.1: Perceptrón de una capa

La arquitectura del perceptrón simple se compone de los siguientes elementos principales:

- **Vector de Pesos:** Un conjunto de parámetros numéricos, denotado por $\vec{w} = (w_1, \dots, w_n) \in \mathbb{R}^n$, que representa la importancia o influencia de cada entrada en la decisión final del perceptrón.
- **Bias (Sesgo):** Un término adicional, $w_0 \in \mathbb{R}$, que ajusta el umbral a partir del cual el perceptrón activa una salida en lugar de otra, permitiendo una mayor flexibilidad en la clasificación.
- **Función de Activación:** Una función matemática, $f : \mathbb{R} \rightarrow \mathbb{R}$, encargada de transformar la suma ponderada de las entradas y el bias en una salida específica del modelo. Para el perceptrón simple, se utiliza comúnmente una función de activación binaria o escalonada, que divide el espacio de entrada en dos categorías distintas.

Dado un vector de entrada $\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$, el perceptrón simple genera una salida $y \in \mathbb{R}$, que se calcula mediante la siguiente expresión:

$$y = P(\vec{x}) = f \left(w_0 + \sum_{i=1}^n w_i x_i \right). \quad (2.1.1)$$

La función de activación escalonada empleada en el perceptrón se define formalmente como:

$$f_k(x) = \begin{cases} 0 & \text{si } x \leq k, \\ 1 & \text{si } x > k, \end{cases} \quad (2.1.2)$$

Esta función mapea la suma ponderada de las entradas a uno de dos posibles estados, 0 o 1, basado en un umbral k . Este mecanismo permite al perceptrón actuar como un clasificador binario, donde la condición de separación entre las dos categorías es determinada por la ecuación del hiperplano $w_0 + \sum_{i=1}^n w_i x_i = k$ en el espacio de entrada \mathbb{R}^n .

Geoméricamente, el perceptrón actúa como un hiperplano en el espacio multidimensional de las entradas, que divide este espacio en dos semiespacios. Cada semiespacio corresponde a uno de los estados posibles (0 o 1). La eficacia del perceptrón depende de que los datos sean linealmente separables, es decir, que exista un hiperplano que pueda dividir sin errores las categorías de datos que intenta clasificar.

2.2. Perceptrón multicapa

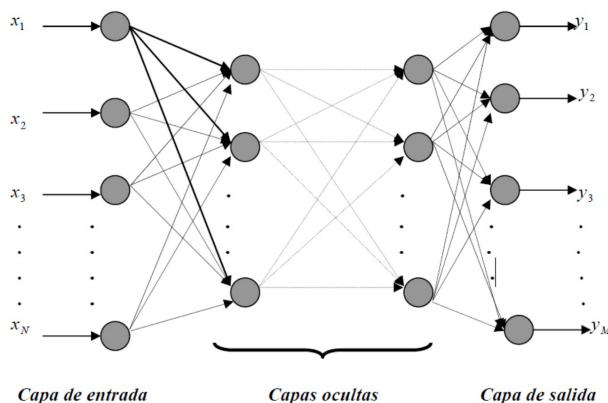


Figura 2.2.1: Perceptrón multicapa

El modelo de perceptrón multicapa (MLP) extiende la arquitectura del perceptrón simple mediante la incorporación de múltiples capas ocultas, lo cual introduce una mayor complejidad tanto en su estructura como en su proceso de entrenamiento. A diferencia del perceptrón simple, cuyo entrenamiento y estructura son relativamente directos al no emplear algoritmos avanzados de optimización, el entrenamiento del MLP se realiza mediante el algoritmo de retropropagación [C].

El proceso de operación del MLP se descompone en dos fases principales: la fase de propagación hacia adelante y la fase de retropropagación del error. En la fase de propagación hacia adelante, la señal de entrada se procesa secuencialmente a través de las capas de la red, empleando funciones de activación específicas en cada nodo, hasta generar la salida en la capa final. Posteriormente, en la fase de retropropagación, el error calculado

2.3. RED CONVOLUCIONAL

como la discrepancia entre la salida obtenida y el valor objetivo se propaga de manera inversa a través de la red, permitiendo la ajuste de los pesos sinápticos.

La estructura del MLP se define formalmente por dos conjuntos principales:

1. Un conjunto de pesos sinápticos,

$$\Theta = \{w_{ij}^l \mid i \in \{1, \dots, n_l\}, j \in \{1, \dots, n_{l-1}\}, l \in \{1, \dots, m\}\},$$

donde w_{ij}^l representa el peso entre el j -ésimo nodo en la capa $l-1$ y el i -ésimo nodo en la capa l .

2. Un conjunto de funciones de activación,

$$F = \{f_i^l : \mathbb{R} \rightarrow \mathbb{R} \mid i \in \{1, \dots, n_l\}, l \in \{1, \dots, m\}\},$$

donde f_i^l designa la función de activación aplicada en el nodo i -ésimo de la capa l .

Aquí, n_l denota el número de nodos en la capa l y m representa el total de capas en la red.

Al igual que en el perceptrón simple, para ejecutar la red es necesario suministrar un conjunto de entradas, $\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$, a partir del cual el MLP produce una salida $y \in \mathbb{R}$. No obstante, la mayor complejidad estructural del MLP implica que su proceso de entrenamiento, fundamentado en el algoritmo de retropropagación, es sustancialmente más elaborado que el correspondiente al perceptrón de una sola capa, destacando la necesidad de una explicación detallada de este procedimiento en la siguiente sección.

2.3. Red convolucional

Las redes neuronales convolucionales [10], son una clase especializada de redes neuronales profundas diseñadas específicamente para procesar datos con una estructura matricial, como imágenes. Han demostrado ser altamente efectivas en una amplia gama de tareas, desde reconocimiento de objetos hasta segmentación de imágenes y más.

2.3.1. Capa convolucional

El corazón de las CNN es la operación de convolución. La convolución implica deslizar un pequeño filtro (o kernel) sobre la imagen de entrada y calcular el producto punto entre los valores de los píxeles de la imagen y los valores del filtro en cada posición. Matemáticamente, la convolución 2D de una imagen I con un filtro K se define como:

$$(I * K)(i, j) = \sum_m \sum_n I(i - m, j - n) \cdot K(m, n)$$

Esta operación produce un mapa de características que resalta ciertas características de la imagen, como bordes, texturas o patrones específicos.

Filtros

Los filtros en una CNN son matrices pequeñas que se aplican a la imagen de entrada. Cada filtro se especializa en detectar ciertas características, como bordes horizontales, verticales o diagonales, colores específicos, texturas, etc. Durante el entrenamiento, la red aprende automáticamente los filtros que son más útiles para cada tarea.

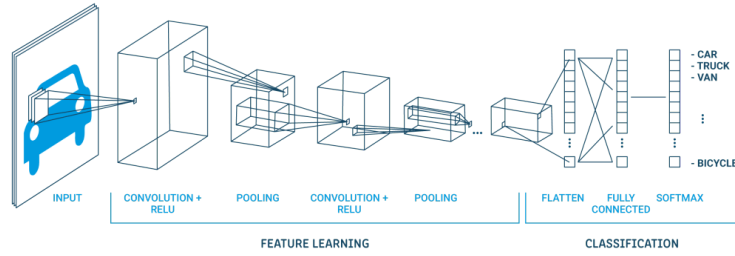


Figura 2.3.1: Estructura de red convolucional

2.3.2. Capas de Pooling

Después de la convolución, es común aplicar una capa de pooling para reducir la dimensionalidad de los datos y hacer que la red sea más robusta ante pequeñas variaciones en la posición de las características detectadas.

Pooling Máximo (Max Pooling)

En el max pooling, se divide la imagen en regiones y se toma el valor máximo de cada región, reduciendo así el tamaño de la representación. Esto se puede expresar como:

$$\text{MaxPooling}(I)(i, j) = \max_{m, n} I(i + m, j + n)$$

Pooling Promedio (Average Pooling)

En el average pooling se toma el valor promedio de cada región, lo que también reduce la dimensionalidad de los datos. Matemáticamente, se define como:

$$\text{AvgPooling}(I)(i, j) = \frac{1}{mn} \sum_m \sum_n I(i + m, j + n)$$

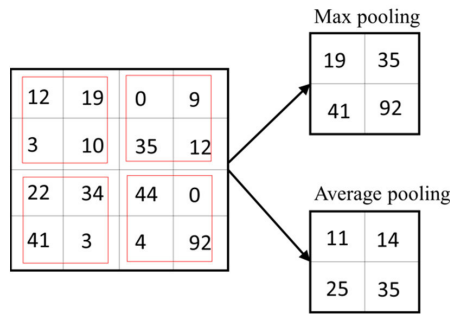


Figura 2.3.2: Pooling

2.4. Red recurrente

Las Redes Neuronales Recurrentes (RNN) son un tipo especializado de redes neuronales diseñadas para modelar secuencias de datos, donde la salida en cada paso de tiempo depende de la entrada actual y de la salida previa. Son ampliamente utilizadas en tareas que involucran datos secuenciales, como el procesamiento del lenguaje natural, la traducción automática y la predicción de series temporales.

La arquitectura básica de una RNN consiste en una capa de neuronas que se conectan consigo mismas a lo largo de una secuencia temporal. Cada neurona en una RNN recibe entradas tanto de la entrada actual como de su estado oculto anterior. Esto permite a la red mantener una memoria interna y capturar dependencias a largo plazo en los datos de secuencia.

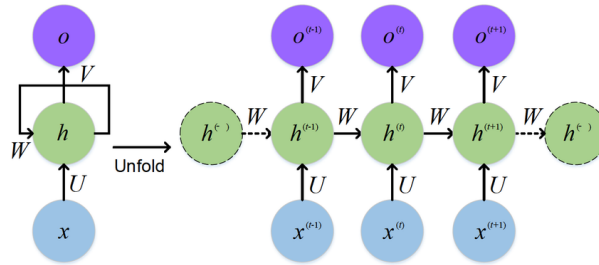


Figura 2.4.1: Estructura de un red neuronal recurrente

2.4.1. Celdas de Memoria en RNN

Celda de Memoria Simple (Simple RNN Cell)

La celda de memoria simple es la forma más básica de una unidad recurrente. En cada paso de tiempo, toma la entrada actual x_t y el estado oculto anterior h_{t-1} , y produce una nueva salida h_t . Matemáticamente, esto se puede expresar como:

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

donde σ es una función de activación no lineal, W_{hx} y W_{hh} son matrices de pesos y b_h es el sesgo.

Celdas de Memoria de Tipo LSTM y GRU

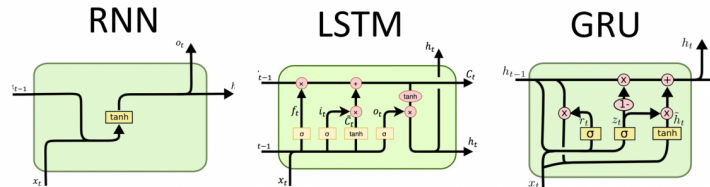


Figura 2.4.2: Distintas estructuras de redes recurrentes

Para abordar el problema del desvanecimiento del gradiente y la dificultad para capturar dependencias a largo plazo, se han propuesto varias variantes de RNN, como Long

Short-Term Memory (LSTM) [16] y Gated Recurrent Unit (GRU) [17]. Estas celdas de memoria están diseñadas para permitir un flujo de información más largo y son capaces de aprender a retener y olvidar información según sea necesario.

Capítulo 3

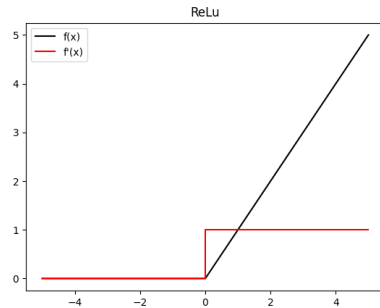
Algunos elementos fundamentales.

3.1. Funciones de activación

Las funciones de activación son responsables de introducir no linealidades en el modelo, por lo que juegan un papel fundamental en las redes neuronales permitiendo que la red aprenda y realice predicciones complejas. La elección de la función de activación adecuada es esencial, dado que impacta directamente en la capacidad de la red para converger durante el entrenamiento y en la precisión de las predicciones que genera.

3.1.1. Unidad Lineal Rectificada (ReLU)

La función ReLU es una función de cálculo rápido que ayuda a mitigar el problema de la saturación de gradiente cuando los valores de la combinación lineal de la entrada son mayores a 0 permitiendo una propagación hacia atrás efectiva incluso en redes profundas.



■ Función:

$$f(x) = \begin{cases} 0 & \text{si } x < 0, \\ x & \text{si } x \geq 0 \end{cases}$$

■ Derivada:

$$f'(x) = \begin{cases} 0 & \text{si } x < 0, \\ 1 & \text{si } x \geq 0 \end{cases}$$

Figura 3.1.1: Relu

3.1.2. SoftMax

Softmax no es una función de activación propiamente dicha, en el sentido que no recibe como entrada un único elemento, sino un vector. Su resultado es un vector cuya suma de valores es igual a 1, simulando así una distribución de probabilidad. La función

y su derivada son más complejas que Relu, necesitando así un mayor coste computacional.

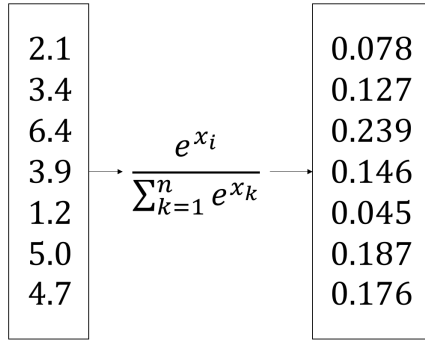


Figura 3.1.2: Softmax

- Función:

$$\alpha(x_j) = \frac{e^{x_j}}{\sum_{k=1}^n e^{x_k}}$$

para $j = 1, \dots, n$.

- Derivada:

$$\alpha'(x_j) = \frac{e^{x_j} (\sum_{k=1}^n e^{x_k}) - e^{2x_j}}{(\sum_{k=1}^n e^{x_k})^2}$$

para $j = 1, \dots, n$.

3.1.3. Sigmoide

Posee la propiedad de mapear el dominio de los reales al intervalo (0,1), con lo cual se asemeja a una distribución de probabilidad. La función sigmoide sufre de un problema conocido como *covariate shift*, debido a que la salida de la función se encuentra centrada en 0.5 existe un desplazamiento de la salida de las neuronas que ocasiona que los pesos tengan problemas de convergencia.

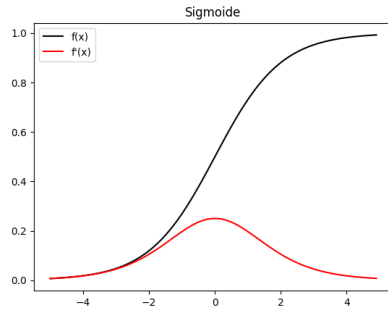


Figura 3.1.3: Sigmoide

- Función:

$$f(x) = \frac{1}{1 + e^{-x}}$$

- Derivada:

$$f'(x) = \frac{e^{-x}}{(1 + e^{-x})^2}$$

3.1.4. Tangente hiperbólica

La tangente hiperbólica se puede ver como una versión escalada y desplazada de la sigmoide. Esta función mapea el dominio de los reales al intervalo abierto (-1, 1). Debido a lo anterior se encuentra centrada en 0 y elimina el problema del *covariate shift*. De igual forma su derivada es más 'activa' que la sigmoide, lo cual permite reducir tiempos de entrenamiento y convergencia, sin embargo, también sufre de saturación de gradiente a los extremos de la función, $\lim_{x \rightarrow \pm\infty} f'(x) = 0$.

3.2. FUNCIÓN DE PÉRDIDA

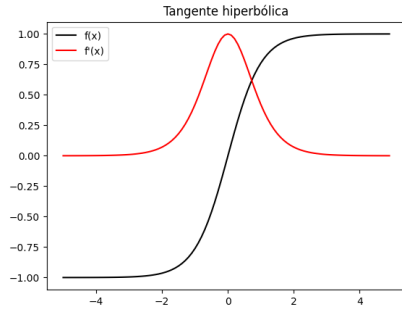


Figura 3.1.4: Tangente Hiperbólica

- Función:

$$f(x) = \frac{2}{1 + e^{-2x}} - 1$$

- Derivada:

$$f'(x) = \frac{4e^{-2x}}{(1 + e^{-2x})^2}$$

3.2. Función de Pérdida

En el aprendizaje automático, especialmente en las redes neuronales, las funciones de pérdida (también conocidas como funciones de costo) juegan un papel crucial. Estas funciones cuantifican el error entre las predicciones realizadas por el modelo y los valores reales observados. El objetivo principal del entrenamiento de una red neuronal es minimizar esta función de pérdida, proceso que guía la optimización de los parámetros del modelo (pesos y sesgos).

La función de pérdida C es una función que toma como entradas el valor predicho por el modelo \hat{y} , y el valor real y , y devuelve un número que representa el coste asociado a la diferencia entre ambos. Matemáticamente, se expresa como:

$$C : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$$

donde \mathcal{Y} es el conjunto de posibles etiquetas o valores reales, y \mathbb{R}^+ denota los números reales no negativos, indicando que el coste no puede ser negativo.

3.2.1. Ejemplos comunes

Error Cuadrático Medio (MSE)

El Error Cuadrático Medio (MSE, por sus siglas en inglés) es ampliamente utilizado en problemas de regresión. Se define como el promedio de los cuadrados de las diferencias entre los valores reales y los predichos:

$$C(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

donde y_i es el valor real y \hat{y}_i es la predicción del modelo para la i -ésima instancia.

Entropía Cruzada

La Entropía Cruzada es una medida de la diferencia entre dos distribuciones de probabilidad, y se utiliza principalmente en tareas de clasificación. Para la clasificación binaria, se define como:

$$C(y, \hat{y}) = - \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

y para la clasificación multiclase, se extiende a:

$$C(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^n \sum_{j=1}^m y_{ij} \log(\hat{y}_{ij})$$

donde y_{ij} indica si la clase j es la clase correcta para la observación i y \hat{y}_{ij} es la probabilidad predicha de que la observación i pertenezca a la clase j .

3.3. Optimizadores

La optimización es una piedra angular en el campo del aprendizaje automático y, en particular, es esencial para el entrenamiento eficaz de redes neuronales. Este proceso matemático se enfoca en encontrar los valores óptimos de los parámetros de un modelo que minimizan o maximizan una función objetivo. En el contexto de las redes neuronales, la función objetivo suele ser la función de pérdida o costo, $L(\Theta)$ [D], que mide el desempeño del modelo al comparar las salidas predichas con los valores reales esperados.

Por lo tanto, el proceso de optimización busca ajustar los parámetros de la red, denotados por Θ , de tal manera que se minimice la función de pérdida. Este enfoque no solo mejora la precisión de las predicciones realizadas por la red, sino que también contribuye a la generalización del modelo, permitiéndole realizar predicciones acertadas sobre datos no vistos anteriormente.

Matemáticamente, el objetivo de la optimización en redes neuronales se puede expresar como la búsqueda del conjunto de valores para Θ que resulte en el mínimo valor de $L(\Theta)$:

$$\min_{\Theta} L(\Theta)$$

donde Θ representa el conjunto de todos los parámetros de la red, que incluyen los pesos y sesgos distribuidos a lo largo de sus diversas capas.

3.3.1. Descenso del Gradiente Estocástico (SGD)

El Descenso del Gradiente Estocástico (SGD) es una variante del método de descenso del gradiente tradicional, optimizado para tratar con grandes volúmenes de datos. A diferencia del descenso del gradiente que requiere el cálculo del gradiente de la función de pérdida sobre el conjunto de datos completo, el SGD actualiza los parámetros utilizando el gradiente basado solo en un subconjunto aleatorio de datos, conocido como mini-lote. Esta metodología reduce significativamente la carga computacional y permite actualizaciones de parámetros más frecuentes, lo que puede acelerar la convergencia en prácticas de entrenamiento de grandes redes neuronales.

Regla de Actualización de Parámetros

$$\Theta_{t+1} = \Theta_t - \eta \nabla_{\Theta} L(\Theta_t, X_{\text{mini-lote}})$$

donde Θ_t representa los parámetros de la red en la iteración t , η es la tasa de aprendizaje, y $\nabla_{\Theta} L(\Theta_t, X_{\text{mini-lote}})$ es el gradiente de la función de pérdida evaluado en un mini-lote $X_{\text{mini-lote}}$.

Cálculo del Gradiente para el mini-lote

Cuando se tiene una función de pérdida compuesta por la suma de varias funciones de pérdida individuales, como en la siguiente expresión:

$$L(\Theta) = \sum_{i=1}^n L(x_i, y_i; \Theta)$$

donde $L(x_i, y_i; \Theta)$ es la función de pérdida asociada a la entrada x_i y la salida deseada y_i , y Θ representa los parámetros (pesos y sesgos) de la red, el gradiente de $L(\Theta)$ con respecto a los parámetros Θ se calcula utilizando la propiedad de linealidad del operador de derivada.

La derivada (o gradiente) de una suma de funciones es igual a la suma de las derivadas (o gradientes) de cada función. Esto se expresa matemáticamente como:

$$\nabla_{\Theta} L(\Theta) = \nabla_{\Theta} \left(\sum_{i=1}^n L(x_i, y_i; \Theta) \right) = \sum_{i=1}^n \nabla_{\Theta} L(x_i, y_i; \Theta)$$

Durante el proceso de backpropagation en el entrenamiento de una red neuronal, se siguen los siguientes pasos:

1. **Propagación hacia adelante:** Cada entrada x_i es procesada por la red para obtener la predicción \hat{y}_i , utilizando los parámetros actuales Θ .
2. **Cálculo del Error:** Se calcula la pérdida $L(x_i, y_i; \Theta)$ para cada par de entrada-salida usando una función de pérdida predefinida.
3. **Propagación hacia atrás (Backpropagation):** Se calcula el gradiente de la pérdida respecto a cada parámetro de la red para cada ejemplo individual, $\nabla_{\Theta} L(x_i, y_i; \Theta)$.
4. **Acumulación de gradientes:** Los gradientes individuales se suman para obtener el gradiente total de la pérdida acumulada:

$$\nabla_{\Theta} L(\Theta) = \sum_{i=1}^n \nabla_{\Theta} L(x_i, y_i; \Theta)$$

5. **Actualización de Parámetros:** Los parámetros de la red se actualizan en dirección opuesta al gradiente acumulado, típicamente utilizando una tasa de aprendizaje η :

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} L(\Theta)$$

Si se desea utilizar el promedio de las pérdidas para regularizar el impacto del tamaño del conjunto de datos, simplemente se divide el gradiente acumulado por el número de ejemplos n , y la actualización de los parámetros se realiza igualmente pero con este gradiente promedio.

3.3.2. Estimación adaptativa del momento (ADAM)

ADAM combina elementos de otros dos optimizadores, RMSprop y Momentum, ajustando las tasas de aprendizaje de los parámetros a través de estimaciones de los primeros y segundos momentos de los gradientes.

Regla de Actualización de Parámetros

$$\Theta_{t+1} = \Theta_t - \frac{\eta_t}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

Donde:

- Θ_t : Vector de parámetros del modelo en la iteración t .
- η_t : Tasa de aprendizaje en el tiempo t , que puede ser ajustada dinámicamente.
- \hat{m}_t : Estimación corregida del primer momento (la media del gradiente) en el tiempo t . Se calcula como:

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

- $m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla L(\Theta_t)$
- \hat{v}_t : Estimación corregida del segundo momento (la varianza no centrada del gradiente) en el tiempo t . Se calcula como:

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

- $v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla L(\Theta_t))^2$
- β_1, β_2 : Factores de decaimiento para los estimados de los momentos. Típicamente, $\beta_1 = 0,9$ y $\beta_2 = 0,999$ son valores recomendados.
- ϵ : Un pequeño número (por ejemplo, 10^{-8}) para evitar divisiones por cero.

Explicación de los Parámetros en ADAM

ADAM ajusta adaptativamente las tasas de aprendizaje de los parámetros a través de estimaciones de los primeros y segundos momentos de los gradientes. Los parámetros clave en ADAM incluyen:

- \hat{m}_t (**Primer momento corregido**) y \hat{v}_t (**Segundo momento corregido**): Estos son utilizados para obtener tasas de aprendizaje adaptativas para cada parámetro. El primer momento \hat{m}_t es una estimación sesgada hacia cero del primer momento (la media del gradiente), y el segundo momento \hat{v}_t es una estimación sesgada hacia cero de los cuadrados del gradiente (varianza no centrada).
- β_1 y β_2 : Son los factores de decaimiento para los estimados de los momentos. β_1 afecta la estimación del primer momento del gradiente, y β_2 maneja la estimación del segundo momento. Típicamente, se utilizan valores de $\beta_1 = 0,9$ y $\beta_2 = 0,999$, proporcionando un buen equilibrio entre estabilidad y reactividad en la actualización de parámetros.
- ϵ : Es un pequeño número añadido al denominador en la actualización de parámetros para evitar la división por cero, típicamente alrededor de 10^{-8} . Este término ayuda a mantener la estabilidad numérica del algoritmo.

3.3. OPTIMIZADORES

Estos componentes de ADAM son críticos para su efectividad y para la adaptación de las tasas de aprendizaje en base al comportamiento de los gradientes, facilitando el entrenamiento eficiente de modelos complejos en aprendizaje profundo.

Capítulo 4

Ejemplos de redes.

En este capítulo presentaremos diversos ejemplos implementados íntegramente en Python. A lo largo de las diferentes secciones mostraremos fragmentos de código ilustrativos para explicar detalles específicos. Para acceder al código completo ir a los códigos del anexo.

4.1. Perceptrón simple - Separación lineal

En este ejemplo [F] vamos a realizar una separación lineal utilizando la arquitectura del perceptrón simple.

```
1 class Perceptron:
2
3     def __init__(self, n_inputs, lr):
4         self.n = n_inputs
5         self.lr = lr
6         self.w = np.random.randn(n_inputs)
7         self.w0 = np.random.randn(1)
8
9     def forward(self, X):
10         pred = np.dot(self.w, X) + self.w0
11         return pred
```

Listing 4.1: Creación del perceptrón

El objetivo es dividir el espacio bidimensional en dos partes, $A = \{r(x) > 0\}$ y $B = \{r(x) \leq 0\}$. Para este ejemplo, utilizaremos la recta $r(x) = mx + n$ en el espacio, con los parámetros $m = -2$ y $n = +12$.

Una vez dividido el espacio en dos, generaremos puntos aleatorios para entrenar el modelo. El perceptrón recibirá dos entradas por cada ejemplo, las coordenadas x e y , y tendrá una salida z que busquemos que sea igual a la salida de la recta $r(x)$. Para entrenarlo, etiquetaremos los puntos del subespacio A con 1 y los del B con -1 , con la intención de que la salida del perceptrón $z = w_1 \cdot x + w_2 \cdot y + w_0$, termine siendo tras el entrenamiento $z = mx - y + n$. Para tener los puntos de A con valores positivos y los de B con valores negativos, utilizaremos una función de activación simple

$$f(x) = \begin{cases} -1 & \text{si } x \leq 0 \\ 1 & \text{si } x > 0 \end{cases} \quad (4.1.1)$$

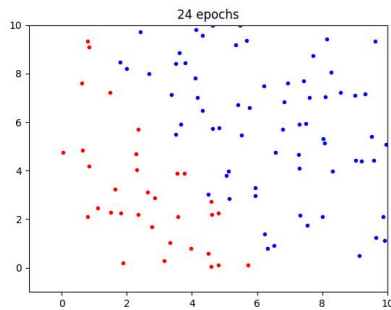
4.1. PERCEPTRÓN SIMPLE - SEPARACIÓN LINEAL

para llevar los resultados a un rango entre -1 y 1.

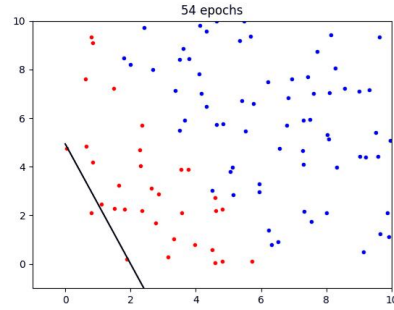
```
1 # functions
2 f = lambda x : -2*x + 12
3 g = lambda x : 1 if x > 0 else -1
4
5 # points
6 px = 10*np.random.rand(n)
7 py = 10*np.random.rand(n)
8
9 points1 = np.array([(x,y) for x,y in zip(px,py) if y > f(x)])
10 points2 = np.array([(x,y) for x,y in zip(px,py) if y <= f(x)])
11
12 # training data
13 X = np.array([px, py])
14 Y = np.array(list(map(g, (py-np.array(list(map(f, px)))))))
```

Listing 4.2: Creación de puntos en el plano

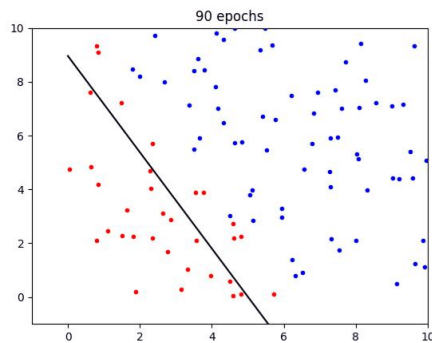
Realizaremos el entrenamiento y comprobaremos cómo el perceptrón separa prácticamente perfectamente el espacio según los puntos de ejemplo que tenemos. Para ello definimos la función *fit* dentro del perceptrón.



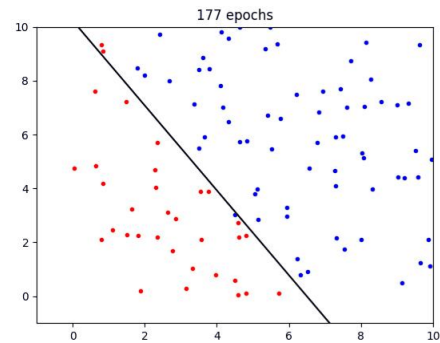
(a) Frame 24



(b) Frame 54



(c) Frame 90



(d) Frame 177

Figura 4.1.1: Evolución del perceptrón durante el entrenamiento

```
1 def fit(self, X, Y, epochs=1):
```

```

2     losses = []
3     for epoch in range(epochs):
4         for i in range(len(X[0])):
5             x = X[:,i]
6             y = Y[i]
7             pred = self.forward(x)
8             error = y-pred
9             self.w += self.lr * error * x
10            self.w0 += self.lr * error
11            losses.append(abs(error))
12            print(f"[{epoch+1}] last error: {losses[-1]} & mean error: {sum(
losses)/len(losses)}")
13    return losses

```

Listing 4.3: Función para entrenar el perceptrón con los puntos del plano

Tras 200 etapas observamos cómo el ajuste del perceptrón es casi perfecto teniendo en cuenta los datos de entrenamiento conocidos.

4.2. Redes Convolucionales - datos MNIST

En este ejemplo [G] exploraremos el uso de redes convolucionales para la clasificación de dígitos escritos a mano utilizando el conjunto de datos MNIST. El conjunto de datos MNIST [22] es una colección de imágenes en escala de grises de dígitos escritos a mano, que van desde 0 hasta 9. Cada imagen tiene una resolución de 28x28 píxeles. Para ello, en python, utilizaremos el dataset proporcionado por keras.

```

1 from keras.datasets import mnist
2 (X, Y), (test_X, test_y) = mnist.load_data()

```

El objetivo principal es familiarizarnos con el uso de redes convolucionales para tareas de clasificación de imágenes. Las redes convolucionales son especialmente efectivas en tareas de visión por computadora debido a su capacidad para aprender y detectar características locales en las imágenes, como bordes, texturas y patrones. Para ello haremos uso de la librería de redes neuronales de pytorch que nos proporcionará la arquitectura necesaria para crear el modelo.

```

1 from torch import nn

```

Utilizaremos una arquitectura de red convolucional relativamente simple que consta de varias capas convolucionales y capas de normalización. Las capas convolucionales son responsables de extraer características relevantes de las imágenes, mientras que las capas de normalización ayudan a estabilizar y acelerar el proceso de entrenamiento.

La arquitectura es la siguiente:

```

1 class Model(nn.Module):
2     def __init__(self, lr):
3         super(Model, self).__init__()
4         self.layer1 = nn.Conv2d(1, 32, 3, padding=1)
5         self.layer2 = nn.BatchNorm2d(32)

```

```
6     self.layer3 = nn.ReLU()
7     self.layer4 = nn.Conv2d(32, 32, 3, padding=1)
8     self.layer5 = nn.BatchNorm2d(32)
9     self.layer6 = nn.ReLU()
10    self.layer7 = nn.Conv2d(32, 1, 3, padding=1)
11    self.layer8 = nn.BatchNorm2d(1)
12    self.layer9 = nn.ReLU()
13    self.layer10 = nn.Flatten()
14    self.layer11 = nn.Linear(28*28, 10)
```

Esta red consta de varias capas convolucionales intercaladas con capas de normalización y activación ReLU. La capa final es una capa lineal que produce la salida final de la red, que representa las probabilidades de cada clase (dígito del 0 al 9).

Tras el entrenamiento obtenemos una precisión del 96.54 %, muy efectiva teniendo en cuenta el tamaño del modelo.

Explorando los Filtros Convolucionales

Una parte interesante de las redes convolucionales es su capacidad para aprender y visualizar los filtros convolucionales que se aplican a las imágenes de entrada. Estos filtros actúan como detectores de características y son esenciales para la capacidad de la red para reconocer patrones en las imágenes.

En la siguiente imagen observamos los 32 filtros convolucionales 3x3 aprendidos por la primera capa de la red.

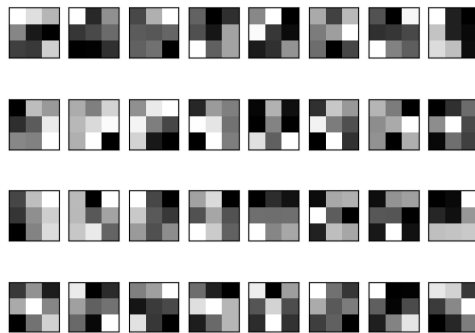


Figura 4.2.1: Filtros de la primera capa de la red

Estos filtros capturan características simples, como bordes, texturas y gradientes, que son fundamentales para el proceso de clasificación. A continuación observamos los resultados de un ejemplo

pasado por los 32 filtros de la red, para observar en qué se está fijando internamente la red. Tenemos dos visualizaciones dependiendo de la escala de color para apreciar mejor cada filtro.

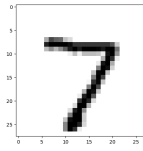


Figura 4.2.2: Ejemplo MNIST

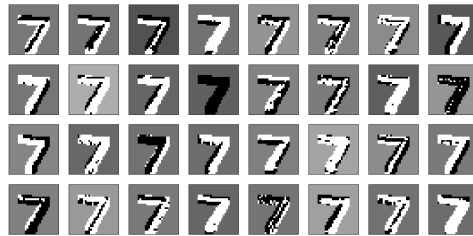


Figura 4.2.3: Ejemplo MNIST tras el primer filtro

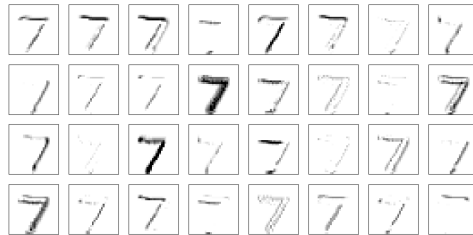


Figura 4.2.4: Ejemplo MNIST tras el primer filtro (V2)

4.3. Redes Recurrentes para Predicción de Series Temporales

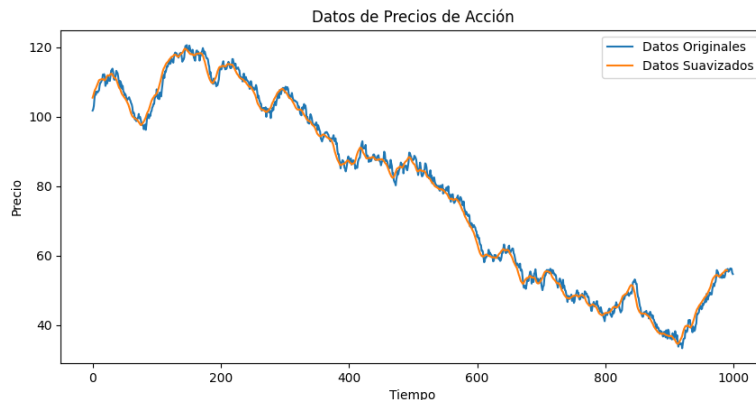
Las redes recurrentes, como las LSTM (Long Short-Term Memory), son ampliamente utilizadas en la predicción de series temporales debido a su capacidad para capturar dependencias temporales en los datos. Estas redes son capaces de recordar información de estados anteriores a medida que procesan nuevas entradas, lo que las hace ideales para modelar patrones y tendencias en series temporales, como los precios de acciones.

En este ejemplo [H], utilizaremos una red LSTM para predecir el precio de una acción en la siguiente hora en función de los precios históricos de la acción. Utilizaremos de nuevo la biblioteca Keras, que proporciona una interfaz sencilla para construir y entrenar redes neuronales.

Primero, generaremos datos de ejemplo que simulan los precios históricos de una acción. Utilizaremos la función `numpy.random.randn()` para generar precios aleatorios y luego aplicaremos una función de suavizado para que los datos se asemejen más a los precios reales de una acción.

4.3. REDES RECURRENTE PARA PREDICCIÓN DE SERIES TEMPORALES

```
1 import numpy as np
2
3 # Generar datos de precios de ejemplo
4 np.random.seed(0)
5 num_datos = 1000
6 precios = 100 + np.cumsum(np.random.randn(num_datos))
7
8 # Funci n de suavizado para simular datos de precios de acciones
9 def suavizar_datos(datos, ventana=10):
10     suavizado = np.convolve(datos, np.ones(ventana)/ventana, mode='valid')
11     return suavizado
12
13 precios_suavizados = suavizar_datos(precios)
14
15 # Visualizar los datos de precios de ejemplo
16 import matplotlib.pyplot as plt
17
18 plt.plot(precios, label='Datos Originales')
19 plt.plot(range(len(precios_suavizados)), precios_suavizados, label='Datos
    Suavizados')
20 plt.xlabel('Tiempo')
21 plt.ylabel('Precio')
22 plt.legend()
23 plt.title('Datos de Precios de Acci n')
24 plt.show()
```



Ahora que tenemos nuestros datos de precios de ejemplo, podemos dividirlos en secuencias de longitud fija para entrenar nuestra red LSTM. Utilizaremos, por ejemplo, las 24 horas anteriores para predecir el precio de la siguiente hora.

```
1 # Funci n para crear secuencias de datos
2 def crear_secuencias(datos, longitud_secuencia):
3     secuencias = []
4     for i in range(len(datos) - longitud_secuencia):
5         secuencia = datos[i:i+longitud_secuencia]
6         objetivo = datos[i+longitud_secuencia]
7         secuencias.append((secuencia, objetivo))
8     return secuencias
9
10 # Definir la longitud de la secuencia
11 longitud_secuencia = 24 # Utilizamos 24 horas anteriores para predecir la
    siguiente
12
```

```

13 # Crear secuencias de datos
14 secuencias = crear_secuencias(precios_suavizados, longitud_secuencia)
15
16 # Dividir los datos en conjuntos de entrenamiento y prueba
17 porcentaje_entrenamiento = 0.8
18 num_entrenamiento = int(porcentaje_entrenamiento * len(secuencias))
19
20 datos_entrenamiento = secuencias[:num_entrenamiento]
21 datos_prueba = secuencias[num_entrenamiento:]
22
23 # Convertir los datos a arreglos numpy
24 X_train, y_train = np.array([x for x, _ in datos_entrenamiento]), np.array(
    ([y for _, y in datos_entrenamiento])
25 X_test, y_test = np.array([x for x, _ in datos_prueba]), np.array([y for _,
    y in datos_prueba])
26
27 # Ajustar las dimensiones de los datos para LSTM (n mero de muestras,
    longitud de la secuencia, n mero de caracter sticas)
28 X_train = np.reshape(X_train, (X_train.shape[0], X_train.shape[1], 1))
29 X_test = np.reshape(X_test, (X_test.shape[0], X_test.shape[1], 1))

```

Con los datos preparados, ahora podemos construir y entrenar nuestra red LSTM utilizando Keras.

```

1 from keras.models import Sequential
2 from keras.layers import LSTM, Dense
3
4 # Definir el modelo LSTM
5 modelo = Sequential()
6 modelo.add(LSTM(50, input_shape=(longitud_secuencia, 1)))
7 modelo.add(Dense(1))
8
9 # Compilar el modelo
10 modelo.compile(optimizer='adam', loss='mean_squared_error')
11
12 # Entrenar el modelo
13 modelo.fit(X_train, y_train, epochs=100, batch_size=32, validation_data=(
    X_test, y_test), verbose=1)

```

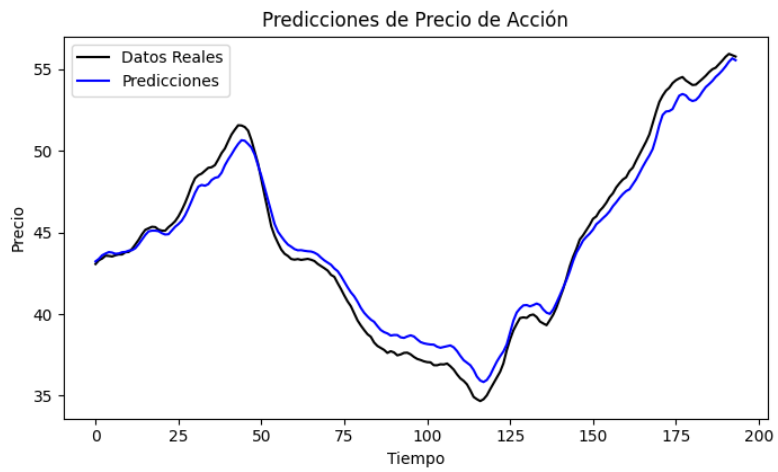
Una vez entrenado el modelo, podemos hacer predicciones sobre los datos de prueba y visualizar los resultados.

```

1 # Hacer predicciones
2 predicciones = modelo.predict(X_test)
3
4 # Visualizar los resultados
5 plt.plot(y_test, label='Datos Reales')
6 plt.plot(predicciones, label='Predicciones')
7 plt.xlabel('Tiempo')
8 plt.ylabel('Precio')
9 plt.legend()
10 plt.title('Predicciones de Precio de Acci n')
11 plt.show()

```


4.3. REDES RECURRENTE PARA PREDICCIÓN DE SERIES TEMPORALES



Parte II

PARTE 2: Aprendizaje por transferencia.

Capítulo 5

Introducción al aprendizaje por transferencia

5.1. Introducción

El aprendizaje por transferencia [20] es una técnica en el campo del aprendizaje automático, especialmente relevante para las redes neuronales profundas. Esta metodología surgió como solución a dos grandes desafíos: la necesidad de extensos volúmenes de datos y el alto coste computacional asociado al entrenamiento de modelos profundos. La premisa del aprendizaje por transferencia consiste en utilizar el conocimiento adquirido (modelos preentrenados) en una tarea para facilitar el aprendizaje en tareas relacionadas pero distintas.

Este enfoque se basa en la observación de que ciertas características generales aprendidas en tareas como el reconocimiento visual son transferibles a una amplia variedad de problemas dentro del mismo dominio. Por ejemplo, las características que detectan bordes, texturas y formas en imágenes son útiles en múltiples tareas de procesamiento visual.

5.2. Fundamentos

Desde una perspectiva matemática, consideramos dos dominios de aprendizaje: el dominio fuente D_S donde el modelo ha sido entrenado, y el dominio objetivo D_T donde se aplicará el modelo. Un modelo de red neuronal profunda generalmente consta de múltiples capas, donde las capas iniciales funcionan como extractores genéricos de características y las capas profundas se especializan más hacia la tarea original.

El proceso de transferencia típicamente implica:

1. **Selección del Modelo Preentrenado:** Elección de un modelo entrenado previamente en un conjunto de datos amplio.
2. **Reutilización de Características:** Uso de las capas iniciales del modelo preentrenado, que contienen características genéricas útiles para la nueva tarea.
3. **Afinamiento (Fine-Tuning):** Ajuste de las capas profundas que se entrenan nuevamente en un conjunto de datos específico para la nueva tarea, ajustando los pesos a partir de los ya aprendidos.

Para implementar esto técnicamente, es posible “congelar” las capas iniciales del modelo, de modo que sus pesos no se actualicen durante el entrenamiento en el nuevo conjunto de datos. Esto se realiza estableciendo los parámetros de estas capas como no entrenables:

```
1 for layer in model.layers[:n]:  
2     layer.trainable = False
```

Capítulo 6

Ejemplos de aprendizaje por transferencia

En este capítulo, exploraremos ejemplos prácticos de aprendizaje por transferencia, con especial atención en su aplicación y efectividad. El segundo ejemplo, enfocado en entrenamiento de procesamiento de lenguaje natural, que es notablemente exigente desde el punto de vista computacional, ha sido desarrollado utilizando un entorno en una cuenta básica de Google Colab para aprovechar el uso de una GPU. Esto nos ha permitido acelerar significativamente el proceso de entrenamiento. Para acceder al código completo ir a los códigos del anexo.

6.1. Clasificación de Imágenes con ImageNet: Gatos vs Perros

En este ejemplo [1], exploraremos cómo utilizar el conjunto de datos ImageNet [4] y la técnica de transferencia de aprendizaje para entrenar un modelo de clasificación binaria que distinga entre imágenes de gatos y perros. Utilizaremos el modelo MobileNetV2 pre-entrenado en ImageNet como punto de partida y adaptaremos las capas finales para nuestro problema específico.

6.1.1. Configuración y Preparación de Datos

En esta sección, configuraremos el entorno y prepararemos los datos para el entrenamiento del modelo.

Comenzaremos importando las bibliotecas necesarias, para la construcción y entrenamiento del modelo.

```
1 from keras.applications import MobileNetV2
2 from keras.models import Sequential
3 from keras.layers import Dense, GlobalAveragePooling2D
4 from keras.preprocessing import image
5 from tensorflow.keras.preprocessing.image import ImageDataGenerator
6 from keras.applications.mobilenet_v2 import preprocess_input,
   decode_predictions
7 import numpy as np
8 import os
```

Después nos descargaremos una base de datos de perros y gatos para poder realizar el entrenamiento de nuestro modelo. Para ello nos iremos a Kaggle y descargaremos el dataset de *Cat VS Dog Dataset* el cual contiene 12491 imágenes de gatos y 12470 imágenes de perros. Una vez descargado lo separaremos en dos datasets, uno para el entrenamiento en el que dejaremos 10.000 imágenes de cada tipo y el resto para la validación, es decir un split de 80 / 20. El sistema de archivos debería quedar de la siguiente forma:

```
1 dataset/  
2   train/  
3     perro/  
4       imagen1.jpg  
5       imagen2.jpg  
6       ...  
7     gato/  
8       imagen1.jpg  
9       imagen2.jpg  
10      ...  
11   validation/  
12     perro/  
13       imagen1.jpg  
14       imagen2.jpg  
15       ...  
16     gato/  
17       imagen1.jpg  
18       imagen2.jpg  
19       ...
```

6.1.2. Preparación del Modelo y Aprendizaje por Transferencia

En esta sección, configuraremos el modelo de clasificación utilizando la técnica de transferencia de aprendizaje y prepararemos las capas finales para la clasificación binaria de gatos y perros.

Utilizaremos el modelo pre-entrenado MobileNetV2, que ha sido entrenado en el conjunto de datos ImageNet para tareas de clasificación de imágenes.

```
1 base_model = MobileNetV2(weights='imagenet', include_top=False, input_shape  
    =(224, 224, 3)) # Excluir las capas finales y especificar el tamaño  
    de entrada
```

y congelaremos todas las capas del modelo pre-entrenado, excepto las capas finales.

```
1 for layer in base_model.layers:  
2     layer.trainable = False
```

Posteriormente crearemos el modelo final añadiendo nuevas capas para adaptar el modelo a nuestra tarea de clasificación de gatos y perros.

```
1 model = Sequential([  
2     base_model,  
3     GlobalAveragePooling2D(),  
4     Dense(128, activation='relu'),
```

6.1. CLASIFICACIÓN DE IMÁGENES CON IMAGENET: GATOS VS PERROS

```
5 Dense(1, activation='sigmoid') # Capa final con activaci n sigmoide
6 ])
```

Para obtener un resumen del modelo podemos hacer *model.summary()* y obtenemos lo siguiente:

```
1 Model: "sequential"
2
3 Layer (type)                Output Shape
4 Param #
5 mobilenetv2_1.00_224 (Functional) (None, 7, 7, 1280)
6 2,257,984
7 global_average_pooling2d    (None, 1280)
8 0
9 (GlobalAveragePooling2D)
10 dense (Dense)               (None, 128)
11 163,968
12 dense_1 (Dense)             (None, 1)
13 129
14 Total params: 2,750,277 (10.49 MB)
15 Trainable params: 164,097 (641.00 KB)
16 Non-trainable params: 2,257,984 (8.61 MB)
17 Optimizer params: 328,196 (1.25 MB)
```

6.1.3. Entrenamiento del Modelo

En esta sección, entrenaremos el modelo utilizando los datos preparados y evaluaremos su rendimiento en un conjunto de datos de validación.

Primero compilaremos el modelo utilizando una función de pérdida binaria y el optimizador adam.

```
1 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['
  accuracy'])
```

Finalmente creamos la función de entrenamiento. En ella primero generamos los objetos *train_generator* y *validation_generator* que contienen las imágenes de nuestro dataset más variaciones de las mismas con giros, ampliaciones, etc, un práctica muy común a la hora de trabajar con imágenes para poder duplicar y triplicar el tamaño de nuestros datasets de manera inteligente. Una vez creados dichos objetos entrenamos el modelo con la función implícita de keras con *model.fit()*.

```

1 def entrenar_modelo(modelo, directorio_entrenamiento, directorio_validacion
  , epochs=10, batch_size=32):
2     # Configurar generadores de datos para entrenamiento y validación con
  aumento de datos
3     train_datagen = ImageDataGenerator(
4         rescale=1./255,
5         shear_range=0.2,
6         zoom_range=0.2,
7         horizontal_flip=True
8     )
9
10    test_datagen = ImageDataGenerator(rescale=1./255)
11
12    train_generator = train_datagen.flow_from_directory(
13        directorio_entrenamiento,
14        target_size=(224, 224),
15        batch_size=batch_size,
16        class_mode='binary' # Especificar la clase binaria (0 para gato, 1
  para perro)
17    )
18
19    validation_generator = test_datagen.flow_from_directory(
20        directorio_validacion,
21        target_size=(224, 224),
22        batch_size=batch_size,
23        class_mode='binary' # Especificar la clase binaria (0 para gato, 1
  para perro)
24    )
25
26    # Entrenar el modelo
27    modelo.fit(
28        train_generator,
29        steps_per_epoch=train_generator.samples // batch_size,
30        epochs=epochs,
31        validation_data=validation_generator,
32        validation_steps=validation_generator.samples // batch_size
33    )

```

6.1.4. Evaluación y Conclusión



Figura 6.1.1: Predicciones: Gato vs Perro

Podemos observar cómo tras un entrenamiento básico de 1'5h hemos conseguido crear

una red casi perfecta, respecto al tipo de imágenes de nuestro dataset, que diferencie entre perros y gatos. Esto es gracias al aprovechamiento de la red *imagenet* que lleva por detrás un entrenamiento de millones de imágenes y horas de entrenamiento con un poder computacional mucho superior al de mi ordenador, por lo que sería impensable obtener estos resultados si quisiera entrenar una red desde cero.

6.2. Análisis de sentimientos con BERT

En el dominio del procesamiento del lenguaje natural (NLP), la adaptación y ajuste fino de modelos preentrenados a tareas específicas han demostrado un éxito significativo. Este enfoque, conocido como aprendizaje por transferencia, aprovecha las características aprendidas de un dominio de problemas para resolver problemas similares con menos gasto computacional que entrenar un modelo desde cero. Este capítulo [J] explora la implementación del aprendizaje por transferencia para el análisis de sentimientos utilizando el modelo BERT (Bidirectional Encoder Representations from Transformers), un método de vanguardia para tareas de NLP.

6.2.1. Configuración y Preparación de Datos

El proceso comienza en un entorno de Google Colab, facilitando el acceso a recursos computacionales potentes como las GPUs. El primer paso implica montar Google Drive para acceder al conjunto de datos almacenado en la nube:

```
1 from google.colab import drive
2 drive.mount('/content/drive', force_remount=True)
```

El conjunto de datos utilizado es "SMSSpamCollection", que contiene una mezcla de mensajes de spam y ham (no spam). La carga de datos se ejecuta con pandas, una biblioteca de manipulación de datos:

```
1 import pandas as pd
2 import os
3
4 main_path = "/content/drive/MyDrive/Colab Notebooks/"
5 file = "SMSSpamCollection.txt"
6 df = pd.read_csv(os.path.join(main_path, "data", file), sep='\t', names=["
    category", "text"])
7 df["label"] = df["category"].astype("category").cat.codes
```

Este fragmento carga los datos, asignando códigos categóricos a las etiquetas—'ham' se convierte en 0 y 'spam' en 1, facilitando el entrenamiento del modelo.

6.2.2. Preparación del Modelo y Aprendizaje por Transferencia

El modelo BERT empleado es "distilbert-base-uncased"[19], una versión más ligera de BERT que retiene la mayor parte del rendimiento del modelo original. Este modelo es adecuado para entornos con recursos computacionales limitados, como dispositivos móviles. Antes de utilizar el modelo, se debe cargar el tokenizador específico para BERT:

```

1 from transformers import DistilBertTokenizerFast
2 tokenizer = DistilBertTokenizerFast.from_pretrained('distilbert-base-uncased')

```

La tokenización [12] convierte el texto en un formato compatible con los requisitos de entrada del modelo. Después de la tokenización, el conjunto de datos se divide en conjuntos de entrenamiento, prueba y validación para garantizar que el rendimiento del modelo sea evaluado adecuadamente:

```

1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(df['text'], df['label'],
3                                                    test_size=0.25, random_state=0)
4 X_test, X_val, y_test, y_val = train_test_split(X_test, y_test, test_size=0.50,
5                                                  random_state=0)

```

6.2.3. Entrenamiento del Modelo BERT

Con los datos preparados y el modelo preentrenado cargado, el siguiente paso involucra adaptar BERT a la tarea de análisis de sentimientos. Esto se hace mediante ajuste fino, donde el modelo preentrenado se entrena (o ajusta) más sobre un conjunto de datos específico. Es crucial notar que solo las capas superiores del modelo se entrenan, ya que son más específicas de la tarea, mientras que las capas inferiores permanecen congeladas:

```

1 from transformers import TFDistilBertForSequenceClassification
2 model = TFDistilBertForSequenceClassification.from_pretrained('distilbert-base-uncased', num_labels=2)
3
4 model.layers[0].trainable = False # Congela las capas del transformador
5 model.layers[1].trainable = False
6 model.summary()

```

```

1 Model: "tf_distil_bert_for_sequence_classification"
2 -----
3 Layer (type)                Output Shape                Param #
4 -----
5 distilbert (TFDistilBertMainLayer)  multiple                    66362880
6
7 pre_classifier (Dense)         multiple                    590592
8
9 classifier (Dense)             multiple                    1538
10
11 dropout_19 (Dropout)          multiple                    0
12
13
14 =====
15 Total params: 66,955,010
16 Trainable params: 1,538
17 Non-trainable params: 66,953,472
18 -----

```

Listing 6.1: Resumen del modelo”

6.2. ANÁLISIS DE SENTIMIENTOS CON BERT

Este enfoque reduce el tiempo y el costo computacional de entrenamiento mientras se adapta el modelo para entender efectivamente el contexto de los mensajes de spam versus ham. Posteriormente, el modelo se compila y se entrena usando los conjuntos de datos preparados:

```
1 model.compile(optimizer="adam")
2 model.fit(train_dataset, validation_data=val_dataset, epochs=3, batch_size
    =64)

1 Epoch 1/3
2 4179/4179 [=====] - 1442s 345ms/step - loss:
    0.1396 - val_loss: 0.0739
3 Epoch 2/3
4 4179/4179 [=====] - 1434s 343ms/step - loss:
    0.0734 - val_loss: 0.0596
5 Epoch 3/3
6 4179/4179 [=====] - 1456s 348ms/step - loss:
    0.0645 - val_loss: 0.0587
```

Listing 6.2: Tiempo total de entrenamiento con una cuenta básica de Google Colab:
1h 12m

6.2.4. Evaluación y Conclusión

Una vez entrenado el modelo, se realizan predicciones sobre el conjunto de prueba y se evalúa su desempeño mediante una matriz de confusión:

```
1 from transformers import TextClassificationPipeline
2 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
3 import numpy as np
4
5 pipeline = TextClassificationPipeline(model=model, tokenizer=tokenizer)
6 pred = pipeline(X_test)
7 pred_labels = [int(p["label"][-1]) for p in pred]
8
9 cm = confusion_matrix(y_test, pred_labels)
10 cm_display = ConfusionMatrixDisplay(cm).plot()
```

obteniendo los siguientes resultados

con una precisión superior al 95 %, este modelo demuestra su robustez. Sin embargo, entrenar un modelo de esta envergadura desde cero sería una tarea descomunal, especialmente con recursos computacionales limitados. Por ejemplo, el entrenamiento de solo 1,538 parámetros me ha llevado aproximadamente 1.2 horas. Extrapolando de manera directa, lo cual es meramente ilustrativo ya que las capas más profundas requieren un esfuerzo computacional mayor debido a la complejidad de calcular los gradientes y al manejo de un volumen significativamente mayor de parámetros y variables en memoria, entrenar el modelo completo con 66,955,010 parámetros tomaría cerca de seis años. Esto subraya la magnitud del modelo que estamos utilizando y la inviabilidad de tal empresa con recursos modestos. Además, alcanzar una precisión tan alta con un conjunto de datos limitado sería prácticamente imposible sin recurrir a técnicas avanzadas de aprendizaje

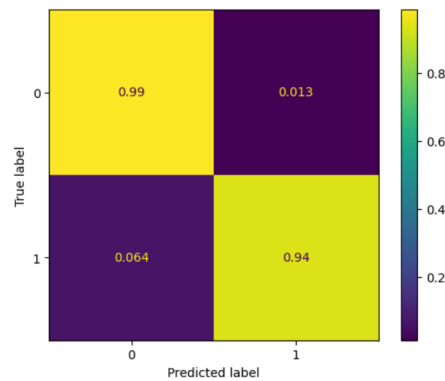


Figura 6.2.1: Matriz de confusión

automático. En este caso específico, estamos utilizando aprendizaje por transferencia con el modelo preentrenado DistilBERT, aplicado al análisis de sentimientos. Estos modelos preentrenados ya poseen una comprensión significativa del lenguaje debido a su entrenamiento previo en grandes corpus de texto, lo que les permite captar matices lingüísticos y contextuales complejos.

Al utilizar un modelo como DistilBERT, podemos ajustar las últimas capas de la red para adaptarlas a nuestras necesidades específicas, en este caso, para diferenciar entre diferentes sentimientos expresados en el texto. Esto nos permite aprovechar el conocimiento preexistente del modelo sobre el lenguaje, ajustando únicamente componentes menores para afinar su rendimiento en tareas específicas. Sin esta base de conocimiento previo, entrenar un modelo desde cero con un conjunto de datos limitado no solo sería impracticable en términos de recursos y tiempo, sino que también correríamos un alto riesgo de sobreajuste (overfitting).

El sobreajuste ocurre cuando un modelo aprende demasiado bien los detalles y el ruido específicos del conjunto de entrenamiento hasta el punto de perder la capacidad de generalizar sobre datos nuevos y no vistos. Esto puede resultar en un modelo que, aunque funcione excepcionalmente bien en los datos de entrenamiento, fallaría al enfrentarse a datos del mundo real o a situaciones fuera de ese conjunto de entrenamiento inicial. Los modelos preentrenados ayudan a mitigar este riesgo al proporcionar una base sólida y generalizada de conocimientos que, con solo pequeños ajustes, se puede adaptar eficazmente a una variedad de tareas sin el peligro inherente de sobreajustar a un conjunto de datos pequeño y potencialmente sesgado.

Parte III

Bibliografía y Anexos

Anexos

C. Backpropagation

El algoritmo de retropropagación es fundamental en el entrenamiento de redes neuronales artificiales y juega un papel crucial en la optimización de sus pesos y sesgos. Desarrollado en la década de 1970 y popularizado por Rumelhart, Hinton y Williams en 1986, este algoritmo ha sido una piedra angular en el éxito de las redes neuronales profundas.

La retropropagación es un método utilizado para calcular eficientemente el gradiente de la función de pérdida de una red neuronal con respecto a cada uno de sus pesos y sesgos. Este gradiente es luego utilizado por algoritmos de optimización, como el descenso del gradiente, para ajustar los parámetros de la red de modo que el error de predicción se minimice.

Denotamos:

- L : número de capas de la red neuronal.
- n_i : número de neuronas en la capa i , donde $i = 1, \dots, L$.
- $x \in \mathbb{R}^{n_1}$: entrada a la red.
- $y \in \mathbb{R}^{n_L}$: salida objetivo (para la entrada x).
- $g : \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_L}$: función que simula la red neuronal.
- C : función de pérdida o función de coste.
- $W^l = (w_{jk}^l)$: los pesos entre capa l y $l + 1$, donde w_{jk}^l es el peso entre el k -ésimo nodo en la capa l y el j -ésimo nodo en la capa $l + 1$.
- f^l : funciones de activación en la capa l .

En la derivación de la retropropagación se utilizan otras cantidades intermedias; se introducen según sea necesario a continuación. Los bias no reciben un tratamiento especial, ya que corresponden a un peso con una entrada fija de 1. A los efectos de la retropropagación, la función de pérdida específica y las funciones de activación no importan, siempre que estas y sus derivadas puedan evaluarse de manera eficiente. Las funciones de activación tradicionales incluyen, entre otras, sigmoid, tanh y ReLU.

La red general es una combinación de composición de funciones y multiplicación de matrices:

$$g(x) = f^L(W^{L-1}f^{L-1}(W^{L-2} \dots f^1(W^1x) \dots)) \quad (\text{C.1})$$

C. BACKPROPAGATION

Para un conjunto de entrenamiento habrá un conjunto de pares de entrada-salida, $\{(x, y) \in \mathbb{R}^{n_1} \times \mathbb{R}^{n_L}\}$. Para cada par de entrada-salida (x, y) en el conjunto de entrenamiento, la pérdida del modelo en ese par es el costo de la diferencia entre la salida proporcionada por la red, $g(x)$ y la salida objetivo y , $C(y, g(x))$. Para ilustrar un posible ejemplo, podemos definir la función de coste con el error cuadrático medio, $C(y, g(x)) := (g(x) - y)^2$.

Pese a parecer que la función de coste depende de la entrada x y la salida deseada y , realmente la variable dependiente en este caso son los pesos ω puesto que queremos ver para qué valores de los pesos se minimiza el error dejando así en cada entrenamiento a x e y como constantes.

Para encontrar el mínimo de la función necesitamos calcular el gradiente de dicha función y por lo tanto encontrar todas las derivadas parciales. Como se puede demostrar / observar, el cálculo individual de todas las derivadas parciales, pese a ser teóricamente posible, necesita de un gasto computacional muy alto. Aquí aparece el algoritmo de Backpropagation en el cual se ayuda de la regla de la cadena reusando cálculos hechos previamente.

La retropropagación se puede expresar para redes simples de retroalimentación en términos de multiplicación de matrices, o más generalmente en términos del gráfico adjunto.

C.1. Algoritmo

En este apartado explicaremos el algoritmo, centrado en un caso básico de una red feedforward, donde los nodos en cada capa están conectados solo a los nodos en la siguiente capa inmediata (sin omitir ninguna capa), y hay una función de pérdida que calcula una pérdida escalar para la salida final, la retropropagación puede ser entendido simplemente por la multiplicación de matrices. Esencialmente, la retropropagación evalúa la expresión de la derivada de la función de costo como un producto de las derivadas entre cada capa de derecha a izquierda, 'hacia atrás', siendo el gradiente de los pesos entre cada capa una simple modificación de los productos parciales.

El objetivo principal de este algoritmo es calcular la derivada parcial de la función de coste respecto a un peso cualquiera de la red, $\frac{\partial C}{\partial \omega_{kj}^l}$. Recordemos primero algunas definiciones:

$$z_k^{l+1} := b_k^{l+1} + \sum_{j=1}^{n_l} \omega_{kj}^l a_j^l \quad (\text{C.2})$$

$$a_k^l := f(z_k^l) \quad (\text{C.3})$$

$$\delta_k^l := \frac{\partial C}{\partial z_k^l} \quad (\text{C.4})$$

$$C := \sum_{i=1}^{n_L} (a_i^L - y_i)^2 = \sum_{i=1}^{n_L} (f(z_i^L) - y_i)^2 \quad (\text{C.5})$$

recordar que la función de coste C no tiene por qué ser la de minimos cuadrados, simplemente la introducimos por el momento para tener una idea más visual. Al ser una función que depende de los pesos, se puede confundir su funcionalidad ya que la notación

completa es $C(\omega_{kj}^l \mid l \in \{1, \dots, L-1\}, j \in \{1, \dots, n_l\}, k \in \{1, \dots, n_{l+1}\})$, la cual puede ser un poco engorrosa.

Tener en cuenta que los sesgos b_k^l también son pesos y tendríamos que calcular sus derivadas parciales, sin embargo pueden calcular igual que los otros simplemente teniendo en cuenta que la entrada de la red será un 1 siempre, por lo que podemos obviar el calculo particular de estas derivadas puesto que no son más que un cálculo particular de los otros.

Aplicando la regla de la cadena tenemos que

$$\frac{\partial C}{\partial \omega_{kj}^l} = \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial \omega_{kj}^l} \stackrel{(C,2)}{=} \frac{\partial C}{\partial z_k^{l+1}} a_j \stackrel{(C,4)}{=} \delta_k^{l+1} a_j \quad (C.6)$$

por lo que para calcular la derivada necesitamos calcular los delta. Veamos el delta de la última capa

$$\begin{aligned} \delta_k^L &:= \frac{\partial C}{\partial z_k^L} \\ &\stackrel{(C,5)}{=} \frac{\partial}{\partial z_k^L} \left(\sum_{i=1}^{n_L} (a_i^L - y_i)^2 \right) \\ &\stackrel{(C,3)}{=} \frac{\partial}{\partial z_k^L} \left(\sum_{i=1}^{n_L} (f(z_i^L) - y_i)^2 \right) \\ &= \sum_{i=1}^{n_L} \frac{\partial}{\partial z_k^L} (f(z_i^L) - y_i)^2 \\ &= \frac{\partial}{\partial z_k^L} (f(z_k^L) - y_k)^2 \\ &= 2(f(z_k^L) - y_k) f'(z_k^L) \in \mathbb{R} \end{aligned} \quad (C.7)$$

que no depende de los pesos. Por otro lado tenemos que

$$z_i^{l+1} = \sum_{s=1}^{n_{l+1}} \omega_{is}^l f(z_s^l) \Rightarrow \frac{\partial z_i^{l+1}}{\partial z_k^l} = \frac{\partial}{\partial z_k^l} \omega_{ik}^l f(z_k^l) = \omega_{ik}^l f'(z_k^l) \quad (C.8)$$

así

$$\begin{aligned}
\delta_k^l &:= \frac{\partial C}{\partial z_k^l} \\
&= \sum_{i=1}^{n_{l+1}} \frac{\partial C}{\partial z_i^{l+1}} \frac{\partial z_i^{l+1}}{\partial z_k^l} \\
&\stackrel{(C,4)}{=} \sum_{i=1}^{n_{l+1}} \delta_i^{l+1} \frac{\partial z_i^{l+1}}{\partial z_k^l} \\
&\stackrel{(C,8)}{=} \sum_{i=1}^{n_{l+1}} [\delta_i^{l+1} \omega_{ik}^l f'(z_k^l)] \\
&= f'(z_k^l) \sum_{i=1}^{n_{l+1}} [\delta_i^{l+1} \omega_{ik}^l]
\end{aligned} \tag{C.9}$$

Ahora bien puesto que $\delta_k^l = h(\delta_k^{l+1})$ para cierta función h definida en (C,9) y $\delta_k^l \in \mathbb{R}$, ya podemos calcular todos los valores de delta con una simple operación. Y con ello directamente por (C,6) podemos calcular todas las derivadas parciales.

D. Aclaración de $L(\Theta)$

Puede surgir la pregunta de por qué he denotado la función de error como $L(\Theta)$ en lugar de $C(y, \hat{y})$, lo cual parecería razonable para mantener la notación de la sección anterior. La función C denota el error cometido entre dos valores reales, típicamente el valor esperado y y el valor predicho \hat{y} . Sin embargo, la función de error en una red neuronal es más compleja que esta simple diferencia.

Para la función de error $L(\Theta)$, asumimos que el valor de salida esperado es fijo, mientras que el valor real de la salida de la red es una expresión completa que depende de los parámetros de la red. Estos parámetros, denotados por Θ , incluyen los pesos y biases de la red. Por lo tanto, $L(\Theta)$ es una función que depende de estos parámetros.

En una red neuronal, los valores de entrada se consideran como variables independientes y los pesos de la red como las variables dependientes que necesitamos ajustar. La función de error $L(\Theta)$ nos permite rastrear cómo los pesos influyen en el rendimiento de la red, facilitando así su modificación para mejorar la precisión del modelo. Mediante el uso de técnicas como la retropropagación [C], podemos actualizar los pesos para minimizar el error y, en consecuencia, mejorar el rendimiento general de la red.

E. Clustering

El *clustering* es una técnica fundamental en el campo del análisis de datos y el aprendizaje automático. Se utiliza para agrupar un conjunto de objetos de manera que los objetos en el mismo grupo (o clúster) sean más similares entre sí que a aquellos en otros grupos. Esta técnica es invaluable en diversas aplicaciones, desde la segmentación de clientes en marketing hasta la detección de patrones en bioinformática.

Matemáticamente, el problema de clustering puede definirse de la siguiente manera: dado un conjunto de n puntos $X = \{x_1, x_2, \dots, x_n\}$ en un espacio métrico \mathbb{R}^d , queremos particionar X en k clústeres $C = \{C_1, C_2, \dots, C_k\}$ tal que la similitud intra-clúster

se maximice y la similitud inter-clúster se minimice. La similitud entre puntos se mide usualmente mediante una función de distancia.

E.1. Algoritmos de Clustering

Existen varios algoritmos para realizar clustering, entre los cuales los más comunes son:

- K-means
- DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Cada uno de estos algoritmos tiene sus propias características y se aplica en diferentes contextos según la naturaleza de los datos y los objetivos del análisis.

E.2. K-means

El algoritmo K-means es uno de los métodos de clustering más utilizados debido a su simplicidad y eficiencia. Dado un número fijo de clústeres k , el objetivo es minimizar la suma de las distancias cuadradas de cada punto a su centroide correspondiente.

Proceso del Algoritmo

1. **Inicialización:** Seleccionar k puntos aleatorios como centroides iniciales.
2. **Asignación:** Asignar cada punto al clúster cuyo centroide sea el más cercano, usando la distancia euclidiana como métrica.
3. **Actualización:** Calcular los nuevos centroides como la media de los puntos asignados a cada clúster.
4. **Convergencia:** Repetir los pasos de asignación y actualización hasta que los centroides no cambien significativamente entre iteraciones.

Función Objetivo

La función objetivo que K-means minimiza es la suma de las distancias cuadradas dentro de los clústeres:

$$J = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

donde μ_i es el centroide del clúster C_i .

E.3. DBSCAN

DBSCAN es un algoritmo de clustering basado en densidad que puede encontrar clústeres de forma arbitraria y manejar ruido en los datos.

Proceso del Algoritmo

1. **Identificación de Puntos Nucleares:** Un punto se considera un núcleo si al menos MinPts puntos están dentro de su vecindario ϵ .
2. **Expansión de Clústeres:** A partir de un punto núcleo, se agregan iterativamente todos los puntos denso-conectados para formar un clúster.
3. **Clasificación de Ruido:** Los puntos que no pertenecen a ningún clúster se clasifican como ruido.

Parámetros Clave

- ϵ : Radio de vecindad.
- MinPts : Número mínimo de puntos en el vecindario para ser considerado un núcleo.

E.4. Evaluación de Clustering

Evaluar la calidad de un clustering es crucial y puede hacerse mediante diversas métricas, tales como:

- **Índice de Dunn:** Mide la compactación y separación de los clústeres.
- **Coefficiente de Silueta:** Evalúa qué tan similar es un punto a su propio clúster en comparación con otros clústeres.
- **Inercia:** Utilizada en K-means, mide la suma de las distancias cuadradas de los puntos a sus centroides.

Códigos

F. Perceptrón simple - Separación lineal

F.1. main.py

```
1 from perceptron import Perceptron
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import os
5 import cv2
6
7
8 # functions
9 f = lambda x : -2*x + 12
10 g = lambda x : 1 if x > 0 else -1
11
12 def ask_number():
13     try:
14         n = int(input(" > "))
15         return n
16     except:
17         print("Invalid argument! Try again:")
18         ask_number()
19
20 def ask_path():
21     name = input(" > ")
22     if not os.path.exists(f"images/{name}/"):
23         os.makedirs(f"images/{name}/")
24         os.makedirs(f"images/{name}/frames/")
25     return name
26     else:
27         print("Folder already exists! Please chose an other name:")
28         ask_path()
29
30 def main():
31
32     # ask for parameters
33     print("Number of points:")
34     n = ask_number()
35     print("Number of epochs:")
36     epochs = ask_number()
37     print("Step:")
38     step = ask_number()
39     print("Name to save the video file:")
40     name = ask_path()
41
42     # points
43     px = 10*np.random.rand(n)
44     py = 10*np.random.rand(n)
45
```

F. PERCEPTRÓN SIMPLE - SEPARACIÓN LINEAL

```
46 points1 = np.array([(x,y) for x,y in zip(px,py) if y > f(x)])
47 points2 = np.array([(x,y) for x,y in zip(px,py) if y <= f(x)])
48
49 # training data
50 X = np.array([px, py])
51 Y = np.array(list(map(g, (py-np.array(list(map(f, px)))))))
52
53 perceptron = Perceptron(n_inputs=2, lr=0.001)
54 imgs = []
55
56 for epoch in range(0, epochs, step):
57
58     perceptron.fit(X, Y, epochs=step)
59     p = lambda x : -(perceptron.w[0] * x + perceptron.w0) / perceptron.
w[1]
60
61     plt.clf()
62     plt.title(str(epoch) + " epochs")
63     plt.plot(points1[:,0], points1[:,1], "b.")
64     plt.plot(points2[:,0], points2[:,1], "r.")
65     plt.plot([0,9], [p(0),p(9)], "k-")
66     plt.xlim(-1,10)
67     plt.ylim(-1,10)
68
69     path = f"images/{name}/frames/{epoch}.jpg"
70     plt.savefig(path)
71     img = cv2.imread(path)
72     imgs.append(img)
73
74     height, width = img.shape[:2]
75
76     video = cv2.VideoWriter(f'images/{name}/{name}.mp4', cv2.
VideoWriter_fourcc(*'mp4v'), 5, (width,height))
77
78     for img in imgs:
79         video.write(img)
80
81     video.release()
82
83
84
85
86 if __name__ == "__main__":
87     main()
```

F.2. perceptron.py

```
1 import numpy as np
2
3 class Perceptron:
4
5     def __init__(self, n_inputs, lr):
6         self.n = n_inputs
7         self.lr = lr
8         self.w = np.random.randn(n_inputs)
9         self.w0 = np.random.randn(1)
10
11     def forward(self, X):
12         pred = np.dot(self.w, X) + self.w0
13         return pred
14
```

```

15     def fit(self, X, Y, epochs=1):
16         losses = []
17         for epoch in range(epochs):
18             for i in range(len(X[0])):
19                 x = X[:,i]
20                 y = Y[i]
21                 pred = self.forward(x)
22                 error = y-pred
23                 self.w += self.lr * error * x
24                 self.w0 += self.lr * error
25                 losses.append(abs(error))
26             print(f"[{epoch+1}] last error: {losses[-1]} & mean error: {sum(
losses)/len(losses)}")
27         return losses

```

G. Redes Convolucionales - datos MNIST

G.1. main.py

```

1 from create_model import create_model
2 from test_model import test_model
3
4
5
6 def main():
7     print("Chose an option:")
8     print(" 1) Train model")
9     print(" 2) Test model")
10    x = input(" > ")
11
12    if x == "1":
13        create_model()
14    else:
15        test_model()
16
17
18
19
20 if __name__ == "__main__":
21     main()

```

G.2. create_model.py

```

1 from matplotlib import pyplot as plt
2 import torch as T
3 import os
4
5 from nn import Model
6 from files import save_file, load_file
7
8
9
10 def load_and_save_data():
11
12     from keras.datasets import mnist

```

```
13 (X, Y), (test_X, test_y) = mnist.load_data()
14
15 if not os.path.exists("data"):
16     os.mkdir("data")
17
18 save_file(X, "data/train_X")
19 save_file(Y, "data/train_y")
20 save_file(test_X, "data/test_X")
21 save_file(test_y, "data/test_y")
22
23 return X, Y
24
25 def create_model():
26
27     # load data MNIST
28     if not os.path.exists("data/test_X"):
29         X, Y = load_and_save_data()
30     else:
31         X, Y = load_file("data/train_X"), load_file("data/train_y")
32
33     f = lambda n : [int(i==n) for i in range(10)]
34     X = T.tensor(X, dtype=T.float).unsqueeze(1)
35     Y = T.tensor(list(map(f, Y)), dtype=T.float)
36
37     # ask for parameters
38     print("\nPath to save the model:")
39     save_path = input(" > ")
40     print("\nNumber of epochs:")
41     epochs = int(input(" > "))
42     print("\nBatch size:")
43     batch_size = int(input(" > "))
44
45     # create & train model
46     nn = Model(0.001)
47     losses, Ks = nn.fit(X,Y, epochs=epochs, batch_size=batch_size)
48
49     # save results
50     if not os.path.exists(f"models/{save_path}"):
51         os.mkdir(f"models/{save_path}")
52
53     # plot losses
54     plt.subplot(2,1,1)
55     k = len(losses) // 100
56     y = losses[::k]
57     x = range(len(y))
58     plt.xticks([])
59     plt.plot(x, y)
60     plt.title("Training loss")
61
62     plt.subplot(2,1,2)
63     k = len(Ks) // 100
64     y = Ks[::k]
65     x = range(len(y))
66     plt.xticks([])
67     plt.plot(x, y)
68     plt.title("Learning rate modificado (K)")
69
70     plt.savefig(f"models/{save_path}/stats.jpg", dpi=200, bbox_inches='
tight')
71     nn.save_model(f"models/{save_path}/model.pkl")
72
73     # show plot
74     plt.show()
```

G.3. test_model.py

```
1 import torch as T
2 import numpy as np
3
4 from nn import Model
5 from files import load_file
6
7
8 def test_model():
9
10     # load model
11     print("\nModel path:")
12     load_path = input(" > ")
13     nn = Model(0.001)
14     nn.load_model(f"models/{load_path}/model.pkl")
15
16     # load test data
17     X_test, y_test = load_file("data/test_X"), load_file("data/test_y")
18     X_test = T.tensor(X_test, dtype=T.float).unsqueeze(1).unsqueeze(1).to(
19         nn.device)
20
21     # make predictions
22     correct = 0
23     for x, y in zip(X_test, y_test):
24         pred = np.argmax(nn(x)[0].tolist())
25         if pred == y:
26             correct += 1
27
28     print("Number of predictions:", correct, f" ({100*correct/len(X_test)}%)")
```

G.4. nn.py

```
1 import torch as T
2 from torch import nn
3 from torch.nn import functional as F
4 from torch import optim
5
6
7 class Model(nn.Module):
8
9     def __init__(self, lr):
10         super(Model, self).__init__()
11
12         # architecture
13         self.layer1 = nn.Conv2d(1,32,3, padding=1)
14         self.layer2 = nn.BatchNorm2d(32)
15         self.layer3 = nn.ReLU()
16         self.layer4 = nn.Conv2d(32,32,3, padding=1)
17         self.layer5 = nn.BatchNorm2d(32)
18         self.layer6 = nn.ReLU()
19         self.layer7 = nn.Conv2d(32,1,3, padding=1)
20         self.layer8 = nn.BatchNorm2d(1)
21         self.layer9 = nn.ReLU()
22         self.layer10 = nn.Flatten()
23         self.layer11 = nn.Linear(28*28, 10)
24         self.layers = [self.layer1, self.layer2, self.layer3, self.layer4,
25             self.layer5, self.layer6, self.layer7, self.layer8, self.layer9, self.
```



```

layer10, self.layer11]
25
26     # optimizer & loss
27     self.optimizer = optim.Adam(self.parameters(), lr=lr)
28     self.loss      = nn.MSELoss()
29
30     # device
31     self.device = T.device('cuda:0' if T.cuda.is_available() else 'cpu'
)
32     self.to(self.device)
33
34 def forward(self, x):
35     for layer in self.layers:
36         x = layer(x)
37     return F.softmax(x, dim=1)
38
39 def fit(self, X, Y, batch_size, epochs):
40
41     n = len(X)
42     losses = []
43     Ks = []
44
45     c = 0
46     for epoch in range(epochs):
47         for i in range(0,n,batch_size):
48             c += 1
49
50             # get sample
51             X_batch = X[i:i+batch_size] if i+batch_size <= n else X[i:]
52             Y_batch = Y[i:i+batch_size] if i+batch_size <= n else Y[i:]
53             X_batch = X_batch.to(self.device)
54             Y_batch = Y_batch.to(self.device)
55
56             # get predictions
57             Z_batch = self.forward(X_batch)
58
59             # train
60             self.optimizer.zero_grad()
61             loss = self.loss(Z_batch, Y_batch).to(self.device)
62             losses.append(loss)
63             loss.backward()
64             self.optimizer.step()
65
66             # Acceso a los gradientes y calculo de K
67             for param in self.parameters():
68                 grad = param.grad.data.flatten()[0] # Gradiente actual
69                 state = self.optimizer.state[param]
70                 if 'exp_avg' in state and 'exp_avg_sq' in state:
71                     exp_avg = state['exp_avg'].flatten()[0] # m_hat
72                     exp_avg_sq = state['exp_avg_sq'].flatten()[0] #
73
74                 v_hat
75
76                 lr = self.optimizer.param_groups[0]['lr']
77                 beta1, beta2 = self.optimizer.param_groups[0]['
betas']
78
79                 eps = self.optimizer.param_groups[0]['eps']
80
81                 # Calculo de K para este parametro
82                 m_hat = exp_avg / (1 - beta1**(c+1))
83                 v_hat = exp_avg_sq / (1 - beta2**(c+1))
84                 K = (lr / (T.sqrt(v_hat) + eps)) * m_hat / grad
85             else:
86                 K = 0
87             Ks.append(K)
88             break

```

```

86         print(f"[{epoch}] -- loss: {losses[-1]} -- mean loss: {sum(
87             losses)/len(losses)}")
88
89         return T.tensor(losses).tolist(), T.tensor(Ks).tolist()
90
91     def save_model(self, path):
92         T.save(self.state_dict(), path)
93
94     def load_model(self, path):
95         self.load_state_dict(T.load(path))

```

G.5. files.py

```

1 import pickle
2
3
4 def save_file(data, path):
5     file = open(path, "wb")
6     pickle.dump(data, file)
7     file.close()
8
9 def load_file(path):
10     file = open(path, "rb")
11     data = pickle.load(file)
12     file.close()
13     return data

```

G.6. ejemplo_layer1.py

```

1 from nn import Model
2 import torch as T
3 from files import load_file
4 import matplotlib.pyplot as plt
5 import numpy as np
6
7
8
9
10 model = Model(0.001)
11 model.load_model("models/model_50/model.pkl")
12
13 x = load_file("data/test_X")
14 x = T.tensor(np.array([x[0]]), dtype=T.float).to(model.device) # shape:
15     torch.Size([1, 28, 28])
16
17 y = model.layer1(x) # shape: torch.Size([32, 28, 28])
18
19 def see_layer1():
20     for i in range(32):
21         plt.subplot(4,8,i+1)
22         img = model.layer1.weight[i][0].tolist()
23         plt.imshow(img, cmap="gray_r", vmin=-1, vmax=1)
24         plt.xticks([])
25         plt.yticks([])
26     plt.show()

```

```
27 def see_ej():
28     img = x[0].tolist()
29     plt.imshow(img, cmap="gray_r", vmin=0, vmax=255)
30     plt.xticks([])
31     plt.yticks([])
32     plt.show()
33
34 def see_ej_tras_layer1():
35     for i in range(32):
36         plt.subplot(4,8,i+1)
37         img = y[i].tolist()
38         plt.imshow(img, cmap="gray_r", vmin=-1, vmax=1)
39         plt.xticks([])
40         plt.yticks([])
41     plt.show()
42
43 def see_ej_tras_layer1_v2():
44     for i in range(32):
45         plt.subplot(4,8,i+1)
46         img = y[i].tolist()
47         plt.imshow(img, cmap="gray_r", vmin=0, vmax=255)
48         plt.xticks([])
49         plt.yticks([])
50     plt.show()
```

H. Redes Recurrentes para Predicción de Series Temporales

H.1. main.py

```
1 # ----- GENERACIÓN DE DATOS ALEATORIA -----
2
3 import numpy as np
4
5 # Generar datos de precios de ejemplo
6 np.random.seed(0)
7 num_datos = 1000
8 precios = 100 + np.cumsum(np.random.randn(num_datos))
9
10 # Función de suavizado para simular datos de precios de acciones
11 def suavizar_datos(datos, ventana=10):
12     suavizado = np.convolve(datos, np.ones(ventana)/ventana, mode='valid')
13     return suavizado
14
15 precios_suavizados = suavizar_datos(precios)
16
17 # Visualizar los datos de precios de ejemplo
18 import matplotlib.pyplot as plt
19
20 def plot_grafica():
21     plt.plot(precios, label='Datos Originales')
22     plt.plot(precios_suavizados, label='Datos Suavizados')
23     plt.xlabel('Tiempo')
24     plt.ylabel('Precio')
25     plt.legend()
26     plt.title('Datos de Precios de Acción')
27     plt.show()
28
```

```

29 # ----- CREAR SECUENCIAS PARA EL MODELO
30 # -----
31 # Funci n para crear secuencias de datos
32 def crear_secuencias(datos, longitud_secuencia):
33     secuencias = []
34     for i in range(len(datos) - longitud_secuencia):
35         secuencia = datos[i:i+longitud_secuencia]
36         objetivo = datos[i+longitud_secuencia]
37         secuencias.append((secuencia, objetivo))
38     return secuencias
39
40 # Definir la longitud de la secuencia
41 longitud_secuencia = 24 # Utilizamos 24 horas anteriores para predecir la
    siguiente
42
43 # Crear secuencias de datos
44 secuencias = crear_secuencias(precios_suavizados, longitud_secuencia)
45
46 # Dividir los datos en conjuntos de entrenamiento y prueba
47 porcentaje_entrenamiento = 0.8
48 num_entrenamiento = int(porcentaje_entrenamiento * len(secuencias))
49
50 datos_entrenamiento = secuencias[:num_entrenamiento]
51 datos_prueba = secuencias[num_entrenamiento:]
52
53 # Convertir los datos a arreglos numpy
54 X_train, y_train = np.array([x for x, _ in datos_entrenamiento]), np.array
    ([y for _, y in datos_entrenamiento])
55 X_test, y_test = np.array([x for x, _ in datos_prueba]), np.array([y for _,
    y in datos_prueba])
56
57 # Ajustar las dimensiones de los datos para LSTM (n mero de muestras,
    longitud de la secuencia, n mero de caracter sticas)
58 X_train = np.reshape(X_train, (X_train.shape[0], X_train.shape[1], 1))
59 X_test = np.reshape(X_test, (X_test.shape[0], X_test.shape[1], 1))
60
61
62 # ----- CREAR MODELO
63 # -----
64 from keras.models import Sequential
65 from keras.layers import LSTM, Dense
66
67 # Definir el modelo LSTM
68 modelo = Sequential()
69 modelo.add(LSTM(50, input_shape=(longitud_secuencia, 1)))
70 modelo.add(Dense(1))
71
72 # Compilar el modelo
73 modelo.compile(optimizer='adam', loss='mean_squared_error')
74
75 # Entrenar el modelo
76 modelo.fit(X_train, y_train, epochs=100, batch_size=32, validation_data=(
    X_test, y_test), verbose=1)
77
78
79 # ----- OBSERVAR RESULTADOS
80 # -----
81 # Hacer predicciones
82 predicciones = modelo.predict(X_test)
83
84 # Visualizar los resultados
85 plt.plot(y_test, "k-", label='Datos Reales')

```

```
86 plt.plot(predicciones, "r-", label='Predicciones')
87 plt.xlabel('Tiempo')
88 plt.ylabel('Precio')
89 plt.legend()
90 plt.title('Predicciones de Precio de Acci n')
91 plt.show()
```

I. Clasificación de Imágenes con ImageNet: Gatos vs Perros

I.1. main.py

```
1 from keras.applications import MobileNetV2
2 from keras.models import Sequential
3 from keras.layers import Dense, GlobalAveragePooling2D
4 from keras.preprocessing import image
5 from tensorflow.keras.preprocessing.image import ImageDataGenerator
6 from keras.applications.mobilenet_v2 import preprocess_input,
   decode_predictions
7 import numpy as np
8 import os
9
10 # Cargar el modelo MobileNetV2 preentrenado
11 base_model = MobileNetV2(weights='imagenet', include_top=False, input_shape
   =(224, 224, 3)) # Excluir las capas finales y especificar el tama o
   de entrada
12
13 # Congelar todas las capas del modelo preentrenado
14 for layer in base_model.layers:
15     layer.trainable = False
16
17 # A adir capas adicionales para la clasificaci n binaria de gato o perro
18 model = Sequential([
19     base_model,
20     GlobalAveragePooling2D(),
21     Dense(128, activation='relu'),
22     Dense(1, activation='sigmoid') # Capa final con activaci n sigmoide
   para la clasificaci n binaria
23 ])
24
25 # Compilar el modelo
26 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['
   accuracy'])
27
28 def predict(img_path):
29     # Cargar una imagen de archivo
30     img = image.load_img(img_path, target_size=(224, 224)) # Redimensionar
   la imagen al tama o que espera MobileNetV2
31     # Convertir la imagen a un array numpy y preprocesarla
32     x = image.img_to_array(img)
33     x = np.expand_dims(x, axis=0)
34     x = preprocess_input(x)
35     # Realizar la predicci n
36     preds = model.predict(x)
37     return 'Gato' if preds[0][0] < 0.5 else 'Perro' # Si la probabilidad
   predicha es mayor a 0.5, se considera un gato; de lo contrario, un
   perro
38
39
```

```

40 # -----
41 # ----- ENTRENAMIENTO -----
42 # -----
43
44 dataset_path = "D:\TFG\imagenes\dogcat"
45
46 def entrenar_modelo(modelo, directorio_entrenamiento, directorio_validacion
47 , epochs=10, batch_size=32):
48     # Configurar generadores de datos para entrenamiento y validaci n con
49     # aumento de datos
50     train_datagen = ImageDataGenerator(
51         rescale=1./255,
52         shear_range=0.2,
53         zoom_range=0.2,
54         horizontal_flip=True
55     )
56
57     test_datagen = ImageDataGenerator(rescale=1./255)
58
59     train_generator = train_datagen.flow_from_directory(
60         directorio_entrenamiento,
61         target_size=(224, 224),
62         batch_size=batch_size,
63         class_mode='binary' # Especificar la clase binaria (0 para gato, 1
64         para perro)
65     )
66
67     validation_generator = test_datagen.flow_from_directory(
68         directorio_validacion,
69         target_size=(224, 224),
70         batch_size=batch_size,
71         class_mode='binary' # Especificar la clase binaria (0 para gato, 1
72         para perro)
73     )
74
75     # Entrenar el modelo
76     modelo.fit(
77         train_generator,
78         steps_per_epoch=train_generator.samples // batch_size,
79         epochs=epochs,
80         validation_data=validation_generator,
81         validation_steps=validation_generator.samples // batch_size
82     )
83
84 import pickle
85
86 # Guardar los pesos en un archivo utilizando pickle
87 def guardar_pesos(pesos, nombre_archivo):
88     with open(nombre_archivo, 'wb') as f:
89         pickle.dump(pesos, f)
90
91 """
92 entrenar_modelo(model, os.path.join(dataset_path, "train"), os.path.join(
93     dataset_path, "train"))
94 guardar_pesos(model.layers[-2].get_weights(), "pesos_dense_1.pkl")
95 guardar_pesos(model.layers[-1].get_weights(), "pesos_dense_2.pkl")
96 """

```

I.2. test.py

```

1 from keras.applications import MobileNetV2
2 from keras.models import Sequential
3 from keras.layers import Dense, GlobalAveragePooling2D
4 from keras.preprocessing import image
5 from keras.applications.mobilenet_v2 import preprocess_input,
   decode_predictions
6
7 import pickle, os
8 import numpy as np
9
10
11 # Cargar los pesos desde un archivo utilizando pickle
12 def cargar_pesos(nombre_archivo):
13     with open(nombre_archivo, 'rb') as f:
14         pesos = pickle.load(f)
15     return pesos
16
17 # Cargar el modelo MobileNetV2 preentrenado
18 base_model = MobileNetV2(weights='imagenet', include_top=False, input_shape
   =(224, 224, 3)) # Excluir las capas finales y especificar el tama o
   de entrada
19
20 # Congelar todas las capas del modelo preentrenado
21 for layer in base_model.layers:
22     layer.trainable = False
23
24 # A adir capas adicionales para la clasificaci n binaria de gato o perro
25 model = Sequential([
26     base_model,
27     GlobalAveragePooling2D(),
28     Dense(128, activation='relu'),
29     Dense(1, activation='sigmoid') # Capa final con activaci n sigmoide
   para la clasificaci n binaria
30 ])
31
32 model.layers[-2].build(input_shape=(None, 1280))
33 model.layers[-2].set_weights(cargar_pesos("model/pesos_dense_1.pkl"))
34
35 model.layers[-1].build(input_shape=(None, 128))
36 model.layers[-1].set_weights(cargar_pesos("model/pesos_dense_2.pkl"))
37
38
39 def predict(img_path):
40     # Cargar una imagen de archivo
41     img = image.load_img(img_path, target_size=(224, 224)) # Redimensionar
   la imagen al tama o que espera MobileNetV2
42     # Convertir la imagen a un array numpy y preprocesarla
43     x = image.img_to_array(img)
44     x = np.expand_dims(x, axis=0)
45     x = preprocess_input(x)
46     # Realizar la predicci n
47     preds = model.predict(x)
48     return 'Gato' if preds[0][0] < 0.5 else 'Perro' # Si la probabilidad
   predicha es mayor a 0.5, se considera un gato; de lo contrario, un
   perro
49
50 # c = lambda x : os.path.join("D:\TFG\imagenes\dogcat\\test\cat",
   "1"+"0000"+str(x))[-4:]+".jpg")
51 # d = lambda x : os.path.join("D:\TFG\imagenes\dogcat\\test\dog",
   "1"+"0000"+str(x))[-4:]+".jpg")

```

J. Análisis de sentimientos con BERT

J.1. sentimental_analysis.ipynb

Initialize Google Colab

```
In [ ]: from google.colab import drive
drive.mount('/content/drive', force_remount=True)
```

Mounted at /content/drive

```
In [ ]: main_path = "/content/drive/MyDrive/Colab Notebooks/"
file = "SMSSpamCollection.txt"
```

Load data and make columns: *text* and *label*

```
In [ ]: import pandas as pd
import os

df = pd.read_csv(os.path.join(main_path, "data", file), sep='\t',
                 names=["category", "text"])
df["label"] = df["category"].astype("category").cat.codes
df.head()
```

```
Out[ ]:
```

	category	text	label
0	ham	Go until jurong point, crazy.. Available only ...	0
1	ham	Ok lar... Joking wif u oni...	0
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	1
3	ham	U dun say so early hor... U c already then say...	0
4	ham	Nah I don't think he goes to usf, he lives aro...	0

```
In [ ]: df.shape
```

```
Out[ ]: (5572, 3)
```

```
In [ ]: X=list(df['text'])
y=list(df['label'])
```

```
In [ ]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, ran
X_test, X_val, y_test, y_val = train_test_split(X_test, y_test, test_size = 0.5)
```

Load tokenizer and Bert model

```
In [ ]: !pip install transformers
```

```
In [ ]: import transformers
print(transformers.__version__)
```

4.30.2

```
In [ ]: from transformers import DistilBertTokenizerFast
tokenizer = DistilBertTokenizerFast.from_pretrained('distilbert-base-uncased')

In [ ]: from transformers import TFDistilBertForSequenceClassification
model = TFDistilBertForSequenceClassification.from_pretrained('distilbert-base-uncased')
```

Preprocess data

```
In [ ]: train_encodings = tokenizer(X_train, truncation=True, padding=True)
test_encodings = tokenizer(X_test, truncation=True, padding=True)
val_encodings = tokenizer(X_val, truncation=True, padding=True)

In [ ]: import tensorflow as tf

train_dataset = tf.data.Dataset.from_tensor_slices((
    dict(train_encodings),
    y_train
))

test_dataset = tf.data.Dataset.from_tensor_slices((
    dict(test_encodings),
    y_test
))

val_dataset = tf.data.Dataset.from_tensor_slices((
    dict(val_encodings),
    y_val
))
```

Fine-tuned Bert model

```
In [ ]: model.layers[0].trainable = False
model.layers[1].trainable = False

In [ ]: model.summary()
```

Model: "tf_distil_bert_for_sequence_classification"

Layer (type)	Output Shape	Param #
distilbert (TFDistilBertMainLayer)	multiple	66362880
pre_classifier (Dense)	multiple	590592
classifier (Dense)	multiple	1538
dropout_19 (Dropout)	multiple	0

=====
Total params: 66,955,010
Trainable params: 1,538
Non-trainable params: 66,953,472
=====

We are going to train only 1,538 parameters. It's very low knowing the task that is been given, but because the Bert architecture part takes care of understanding the language, maybe we dont need much study just to make a simple binary separation.

The fact that i'm choosing to train only those parameters is also because of the training time.

```
In [ ]: model.compile(optimizer="adam")

In [ ]: model.fit(train_dataset, validation_data=test_dataset, epochs=3, batch_size=64)

Epoch 1/3
4179/4179 [=====] - 1442s 345ms/step - loss: 0.1396 - val_loss: 0.0739
Epoch 2/3
4179/4179 [=====] - 1434s 343ms/step - loss: 0.0734 - val_loss: 0.0596
Epoch 3/3
4179/4179 [=====] - 1456s 348ms/step - loss: 0.0645 - val_loss: 0.0587

Out[ ]: <keras.callbacks.History at 0x7f6fed5afa60>
```

Save tokenizer and Model

```
In [ ]: savepath = "/content/drive/MyDrive/Colab Notebooks/models/sentiment_analysis"

In [ ]: model.save_pretrained(savepath)
tokenizer.save_pretrained(savepath)
```

```
Out[ ]: ('/content/drive/MyDrive/Colab Notebooks/ABRICOT/Telcel/NLP/models/sentiment_an
        alysis/tokenizer_config.json',
        '/content/drive/MyDrive/Colab Notebooks/ABRICOT/Telcel/NLP/models/sentiment_an
        alysis/special_tokens_map.json',
        '/content/drive/MyDrive/Colab Notebooks/ABRICOT/Telcel/NLP/models/sentiment_an
        alysis/vocab.txt',
        '/content/drive/MyDrive/Colab Notebooks/ABRICOT/Telcel/NLP/models/sentiment_an
        alysis/added_tokens.json',
        '/content/drive/MyDrive/Colab Notebooks/ABRICOT/Telcel/NLP/models/sentiment_an
        alysis/tokenizer.json')
```

```
In [ ]:
```

Predictions

```
In [ ]: from transformers import TextClassificationPipeline
        from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
        import numpy as np
```

```
In [ ]: pipeline = TextClassificationPipeline(model=model, tokenizer=tokenizer)
```

```
In [ ]: pred = pipeline(X_test)
```

```
In [ ]: print(pred[0])
        int(pred[0]["label"][-1])
```

```
{'label': 'LABEL_0', 'score': 0.8320313096046448}
```

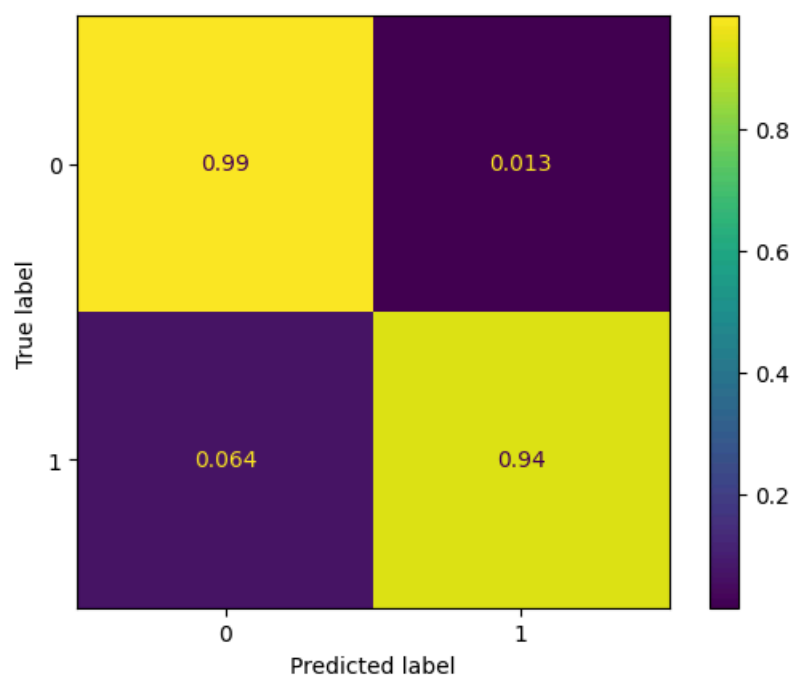
```
Out[ ]: 0
```

```
In [ ]: pred = [int(p["label"][-1]) for p in pred]
```

```
In [ ]: cm = confusion_matrix(y_test, pred)
        cm2 = np.array([cm[i]/cm[i].sum() for i in range(len(cm))])

        cm_display = ConfusionMatrixDisplay(cm2).plot()
        cm
```

```
Out[ ]: array([[594,  8],
               [ 6, 88]])
```



We end up with more than 95% accuracy, so they are good results!

Bibliografía

- [1] Christopher M Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- [2] Jason Brownlee. *Deep Learning for Computer Vision: Image Classification, Object Detection, and Face Recognition in Python*. Machine Learning Mastery, 2019.
- [3] Francois Chollet. *Deep Learning with Python*. Manning Publications, 2017.
- [4] Jia Deng et al. “ImageNet: A large-scale hierarchical image database”. En: *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE. 2009, págs. 248-255.
- [5] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. En: *arXiv preprint arXiv:1810.04805* (2019).
- [6] Xavier Glorot y Yoshua Bengio. “Understanding the difficulty of training deep feed-forward neural networks”. En: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. 2010, págs. 249-256.
- [7] Ian Goodfellow, Yoshua Bengio y Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [8] Simon Haykin. *Neural Networks: A Comprehensive Foundation*. Prentice Hall, 1998.
- [9] John Hertz, Anders Krogh y Richard G Palmer. *Introduction to the Theory of Neural Computation*. Addison-Wesley, 1991.
- [10] Ryan Nash Keiron O’Shea. “An Introduction to Convolutional Neural Networks”. En: *arXiv preprint arXiv:1511.08458* (2015).
- [11] Yu Li y Qiang Yang. “Meta-Learning: A Survey”. En: *arXiv preprint arXiv:1706.04402* (2017).
- [12] Christopher D. Manning e Hinrich Schütze. “Foundations of Statistical Natural Language Processing”. En: MIT Press, 1999. Cap. Tokenization, págs. 124-129.
- [13] Seymour Papert Marvin Minsky. *Perceptrons: an introduction to computational geometry*. 1969.
- [14] Kevin P Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, 2012.
- [15] Sinno Jialin Pan y Qiang Yang. “A Survey on Transfer Learning”. En: *IEEE Transactions on Knowledge and Data Engineering* 22.10 (2010), págs. 1345-1359.
- [16] Eric Rothstein Morris Ralf C. Staudemeyer. “– Understanding LSTM – a tutorial into Long Short-Term Memory Recurrent Neural Networks”. En: *arXiv preprint arXiv:1909.09586* (2019).
- [17] Stefan Güttel Roberto Cahuantzi Xinye Chen. “A comparison of LSTM and GRU networks for learning symbolic sequences”. En: *arXiv preprint arXiv:2107.02248* (2023).
- [18] Sebastian Ruder. “An overview of gradient descent optimization algorithms”. En: *arXiv preprint arXiv:1609.04747* (2016).

BIBLIOGRAFÍA

- [19] Victor Sanh et al. “DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter”. En: *arXiv preprint arXiv:1910.01108* (2019).
- [20] Lisa Torrey y Jude Shavlik. “Transfer Learning”. En: *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques*. IGI Global, 2010, págs. 242-264.
- [21] Ashish Vaswani et al. “Attention is All You Need”. En: *arXiv preprint arXiv:1706.03762* (2017).
- [22] Christopher J.C. Burges Yann LeCun Corinna Cortes. “The MNIST Handwritten Digit Database”. En: *yann.lecun.com* (2020).