



GRADO EN MATEMÁTICA COMPUTACIONAL

TRABAJO FINAL DE GRADO

Red neuronal para la elección de los parámetros de compresión óptimos

Autor:
Marta AICART PAUNER

Supervisor:
Francesc ALTED ABAD
Tutor académico:
Maria Victoria IBAÑEZ GUAL

Fecha de lectura: 20 de junio de 2022
Curso académico 2021/2022

Resumen

Este documento recopila los detalles y procedimientos más importantes del Trabajo Final de Grado, cuya aplicación práctica ha sido desarrollada durante la Estancia en Prácticas en la empresa ironArray.

Uno de los principales objetivos de los científicos a lo largo de la historia ha sido el diseño y construcción de máquinas capaces de realizar procesos con cierta inteligencia. En la actualidad, gracias a muchos años de investigación de profesionales repartidos por todo el mundo, las redes neuronales han alcanzado un gran nivel de madurez y se emplean en todo tipo de aplicaciones.

El objetivo de este Trabajo de Fin de Grado es la implementación de una red neuronal artificial capaz de predecir la mejor combinación de parámetros de compresión para un conjunto de datos dado. Posteriormente, esta red es usada para mejorar el rendimiento del compresor Blosc2 de ironArray. En primer lugar, se proporcionará una base teórica del aprendizaje automático, centrándose en las redes neuronales artificiales. A continuación, se presentará el modelo de red neuronal desarrollado junto con los resultados obtenidos. Después de este análisis, se compararán los objetivos con lo conseguido en la red. Finalmente, se expondrán las conclusiones y líneas futuras de trabajo, a partir de este Trabajo de Fin de Grado.

Palabras clave

Red neuronal, Retropropagación, Gradiente Descendiente.

Keywords

Neural network, Backpropagation, Gradient Descent.

Índice general

1. Introducción	7
1.1. Contexto y motivación del proyecto	7
2. Objetivos	9
3. Desarrollo del TFG	11
3.1. Las redes neuronales	11
3.2. Entrenamiento de las redes neuronales	15
3.3. Categorical Cross Entropy	17
3.4. Adam	17
3.4.1. Gradiente Descendiente (GD)	18
3.4.2. Gradiente Descendiente Estocástico (SGD)	19
3.4.3. Gradiente Descendiente Estocástico con Momentum (SGD con Momentum)	20
3.4.4. Propagación de raíz cuadrática media (RMSProp)	21
3.4.5. Backpropagation	22
3.4.6. Algoritmo Adam	24

3.4.7. Ejemplo	26
4. Resultados	31
4.1. Obtención de los datos	33
4.2. Preparación y manipulación de los datos	35
4.3. Construcción del modelo de red neuronal	37
4.3.1. Arquitectura de la red	37
4.3.2. Compilación de la red	37
4.3.3. Adecuación del modelo (fitting)	37
4.3.4. Resultados de la red	38
4.3.5. Estudio de las observaciones mal clasificadas	40
5. Conclusiones	47
A. Generador de datos	51
A.1. Código entropy_probe.c	51
B. Análisis de datos	67
B.1. Código Neural-Network-CRATIO.ipynb	67
B.2. Código Neural-Network-SPEED.ipynb	81
B.3. Código Neural-Network-BALANCE.ipynb	94

Capítulo 1

Introducción

1.1. Contexto y motivación del proyecto

Este documento recoge el Trabajo Final de Grado, cuyo objetivo principal es el estudio de la elección de los parámetros de compresión óptimos para un conjunto de datos mediante redes neuronales artificiales.

Con el paso del tiempo, para superar limitaciones físicas en el campo de la informática ha sido fundamental el uso de la compresión. Esta técnica minimiza el espacio ocupado por los archivos en disco duro y consecuentemente reduce el tiempo necesario para que sean descargados o transferidos. Esta reducción de espacio y tiempo implica un gran ahorro de costes. Por lo tanto, es fundamental el uso de librerías de rápida compresión como lo es C-Blosc2.

Actualmente, vivimos en un crecimiento tecnológico continuo y, entre estos avances cabe resaltar el logrado en el campo de la inteligencia artificial, concretamente el campo de las redes de neuronas artificiales o, comúnmente llamadas, redes neuronales. Las redes neuronales tratan de imitar el comportamiento de la red neuronal de un cerebro humano. En esta, millones de neuronas forman una red y a través de unas conexiones o sinapsis se comunican entre ellas para realizar distintas acciones. Sin embargo, las redes neuronales artificiales van más allá, ya que tratan de imitar este comportamiento y llevarlo a la computación para automatizar el aprendizaje. Esto supone una gran ventaja ya que no necesitan ser programadas explícitamente y destacan en ámbitos donde la programación convencional tiene recursos limitados para alcanzar soluciones a ciertos problemas.

Este proyecto nace motivado por el interés en comprender el funcionamiento de las redes neuronales y en la oportunidad de poder desarrollar un modelo de red neuronal propio. Este

tema es muy interesante por su gran y veloz desarrollo, y como consecuencia por su uso en una gran variedad de ámbitos tecnológicos. Además, existía la posibilidad de poder aplicar este modelo en una situación real, ya que era idóneo para optimizar el compresor de datos C-Blosc2, una librería que permite comprimir y descomprimir datos a una gran velocidad y sin pérdida de información.

Capítulo 2

Objetivos

El principal objetivo de este trabajo es el estudio y aplicación de las redes neuronales. Con esto se pretende desarrollar un modelo de red neuronal que elija en cada caso la mejor combinación de parámetros para la compresión del meta-compresor C-Blosc2 [8]. Concretamente, dado un dataset se pretende que la red indique cuales son los parámetros de compresión óptimos para maximizar su velocidad y ratio de compresión.

Capítulo 3

Desarrollo del TFG

3.1. Las redes neuronales

La red neuronal es un modelo computacional compuesto por distintos elementos procesales organizados en capas. Estos elementos se denominan neuronas, y a través de canales de comunicación unidireccionales envían información a otras unidades de la red.

La neurona es la unidad básica de procesamiento de la red neuronal. A cada neurona i se le asocia un valor, denominado variable de estado x_i . Tiene conexiones de entrada a través de las cuales recibe los valores de entrada. A estos valores se les aplica una determinada transformación y se genera un valor de salida.

Las neuronas se agrupan por capas, y si estas unidades tienen un conjunto común de entradas se denomina capa densa. Dos neuronas que se encuentren en una misma capa densa reciben la misma información de entrada de la capa anterior, y las salidas de cada capa son los inputs de la capa siguiente. La primera capa se denomina capa de entrada, y es aquella en la que las neuronas reciben los inputs, es decir, los datos a evaluar de cada elemento de la muestra (de cada observación); la última capa, es decir, aquella cuyas neuronas producen el output de la red, se denomina capa de salida; y las capas intermedias se denominan capas ocultas, en las cuales se realizan los cálculos internos de la red. Por lo tanto, la red neuronal consiste en una serie de capas en las que cada una recibe la información procesada por la capa anterior. En la figura 3.1 se puede apreciar la arquitectura descrita de una red neuronal.

Además, cada neurona está conectada con otras a través de unos enlaces, y a cada conexión (i, j) de las neuronas i y j se le asocia un peso $w_{ij} \in \mathbb{R}$. Este peso define con qué intensidad cada variable de entrada afecta a la neurona. De hecho, estos pesos son parámetros del modelo,

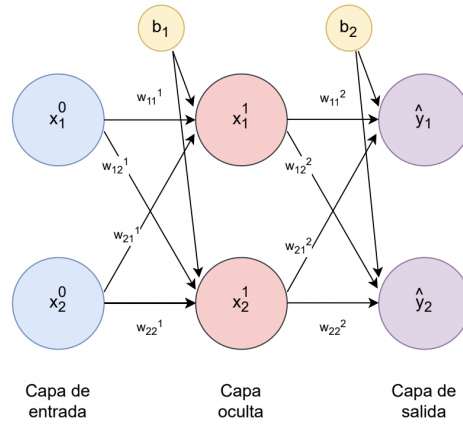


Figura 3.1: Arquitectura de una red neuronal con una única capa oculta

cuyos valores debe estimar la red neuronal.

Estos pesos se representan en forma de matriz, matriz de pesos W , y existe una para cada capa oculta. Su dimensión es $n \times m$ donde, el número de filas n es el número de neuronas de la capa actual y el número de columnas m es el número de neuronas de la capa anterior. Por lo tanto, si W^L es la matriz de pesos de la capa oculta L , entonces el peso w_{ij}^L de esta matriz representa el peso de enlace entre la neurona i de la capa $L - 1$ hacia la neurona j de la capa L .

En la primera capa de la red, la variable de estado x_i de cada neurona es un valor de entrada a la red, y en el resto de capas es el valor obtenido tras realizar una serie de cálculos a partir de los valores de las neuronas de las capas anteriores (Eq. 3.2).

Las variables de estado de las neuronas de la cada capa oculta se puede expresar vectorialmente como $X_L = [x_1^L, x_2^L, \dots, x_m^L]$ para $L = 0, 1, \dots, N$, siendo L la capa oculta, N el número de capas (sin contar la capa de entrada), $L = 0$ la capa de entrada y m el número de neuronas de la capa L , pudiendo variar este número dependiendo de la capa en la que nos situemos.

Para la neurona i de una capa L distinta a la primera, se define la siguiente suma ponderada a partir de las variables de estado de la capa anterior:

$$z_i^L = Z_i^L(X^{L-1}) = \sum_{j=1}^{m_{L-1}} w_{ji}^L x_j^{L-1} + b_L \quad (3.1)$$

donde en cada capa, b_L (bias en inglés) es un nuevo parámetro a estimar. Los sesgos de una

red pueden expresarse en forma vectorial b , donde $b = [b_1, b_2, \dots, b_N]$, siendo b_i el sesgo de la capa oculta i y siendo N el número de capas (sin contar la capa de entrada).

La red conecta múltiples neuronas de forma secuencial, pero en cada neurona se realiza un problema de regresión lineal. Por lo tanto, en la red se concatenan diferentes operaciones de regresión lineal. A la salida de la neurona, puede existir una función no lineal, que modifica el valor resultado o impone un límite que no se debe sobrepasar antes de propagarse a otra neurona. Esta función se conoce como función de activación, que utiliza la suma ponderada de la entrada anterior y la transforma una vez más como salida. Esta función añade deformaciones no lineales al valor de salida de la neurona, para así poder encadenar de forma efectiva la computación de varias unidades. Existen muchas funciones de activación, como por ejemplo la función sigmoide cuya ecuación es $a(z) = \frac{1}{1+e^{-z}}$, la tangente hiperbólica cuya ecuación es $a(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$, etc, pero se va a explicar solo una de ellas, porque es la utilizada en la red neuronal creada que se verá posteriormente. Esta función es la Rectified Linear Unit (ReLU), y se define como

$$x_i^L = a(z_i^L) = \max(0, z_i^L) = \max(0, \sum_{j=1}^{m_{L-1}} w_{ji}^L x_j^{L-1} + b_L) \quad (3.2)$$

Esta función se comporta como una función lineal cuando el valor de entrada es positivo, y constante a cero cuando es negativo. Se trata de una función de activación no lineal, la cual ayuda a modelar funciones curvas o no triviales. En la figura 3.2 se puede apreciar la representación gráfica de diversas funciones de activación.

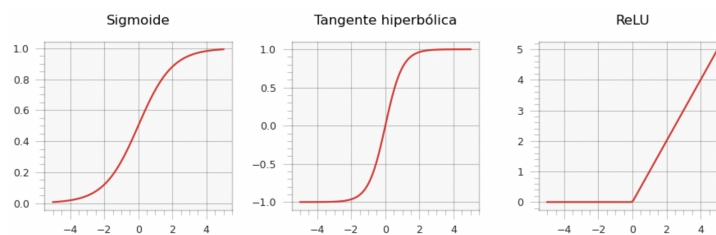


Figura 3.2: Ejemplos de funciones de activación

Por lo tanto, en cada neurona (menos en las neuronas de la capa de entrada) se realiza primero la suma ponderada de sus entradas expresada en la ecuación 3.1, y después a este valor se le aplica una función de activación, es decir, se realiza la composición de dos funciones expresada en la ecuación 3.2. En la figura 3.4 se puede apreciar el proceso de estos cálculos.

Una vez realizados todos los cálculos en las capas ocultas de la red, el resultado de la última capa oculta es utilizado como entrada en la capa de salida para realizar estos mismos cálculos, es decir, los definidos en la ecuación 3.2. En nuestro caso, crearemos una red de clasificación, es decir para clasificar cada observación entre J posibles clases disjuntas. Por lo tanto, la capa

de salida constará de J neuronas, y en cada neurona j se calcula un valor de salida \hat{y}_j , que en nuestro caso nos dará la probabilidad de que la observación que estamos analizando pertenezca a la j -ésima clase. En caso de que $J = 2$, la red tendrá una única neurona en la capa final.

En el caso particular de una red neuronal con 2 neuronas en la capa de entrada, nuestro conjunto de datos será $T = \{(x_1(i), x_2(i), y_i)\}_{i=1, \dots, M}$. Llamaremos muestra u observación a una única fila de datos, la cual contiene unos inputs que son los que se introducen en el algoritmo y una salida, la etiqueta $((x_1(i), x_2(i), y_i), i \in \{1 = 1, \dots, M\})$. Cada observación del conjunto de datos está relacionada con su etiqueta observada, que en la red de clasificación se refiere a la clase a la que pertenece la observación. Tanto a las etiquetas observadas como a los valores predichos por la red se les aplica una codificación denominada *one hot encoding*. Esta codificación consiste en transformar estas variables numéricas en vectores de probabilidad de longitud J , siendo J el número de clases. Estos vectores se expresan como $y_i = [y_{i1}, y_{i2}, \dots, y_{iJ}]$ para las etiquetas observadas, con

$$\sum_{j=1}^J y_{ij} = 1, \quad i = 1, \dots, M, \quad y_{ij} \in \{0, 1\} \quad (3.3)$$

donde todos los elementos del vector son 0, excepto uno, que tiene un 1 como valor, donde el 1 especifica la categoría del elemento.

Utilizaremos la notación $\hat{y}_i = [\hat{y}_{i1}, \hat{y}_{i2}, \dots, \hat{y}_{iJ}]$ para los valores predichos por la red para la muestra i -ésima. Es decir, los valores del vector representan probabilidades y por lo tanto la suma de todos estos valores siempre debe ser 1, es decir:

$$\sum_{j=1}^J \hat{y}_{ij} = 1, \quad i = 1, \dots, M, \quad \hat{y}_{ij} \geq 0 \quad (3.4)$$

Tanto las etiquetas observadas como los valores predichos se representan en matrices con M filas y J columnas, donde M es el número de observaciones, y J es el número de clases. Entonces, el elemento ij -ésimo de las matrices será la probabilidad de que el elemento correspondiente al input i pertenezca a la clase j -ésima. Por lo tanto, el valor y_{ij} es la probabilidad real de que ocurra el evento j -ésimo para la muestra i -ésima y el valor \hat{y}_{ij} es la probabilidad predicha por la red de que ocurra el evento correspondiente para la muestra i -ésima.

Una de las formas de lograr que las salidas de cada capa sean no negativas es utilizar la función de activación ReLU definida anteriormente en la ecuación 3.2. Además, en la última capa utilizaremos la función de activación *softmax* (Ec. 3.5) que toma como entrada un vector de números reales $[x_i]_{i=1}^J$ y lo normaliza en una distribución de probabilidad que consta de J probabilidades proporcionales a las exponenciales de los números de entrada. En este caso, los valores que toma como entrada son los calculados en las neuronas de la capa de salida N , definidos en forma vectorial como $X^N = [x_{i1}^N, x_{i2}^N, \dots, x_{iJ}^N]$.

Por lo tanto, softmax convierte estas salidas en probabilidades, siendo $s(x_{ij}^N)$ la probabilidad de que la muestra i pertenezca a la clase j . Su fórmula es:

$$\hat{y}_{ij} = s(x_{ij}^N) = \frac{e^{x_{ij}^N}}{\sum_{j=1}^J e^{x_{ij}^N}}, \quad j = 1, \dots, J \quad (3.5)$$

donde J es el número de clases. En la figura 3.3 se puede observar el funcionamiento de esta función.

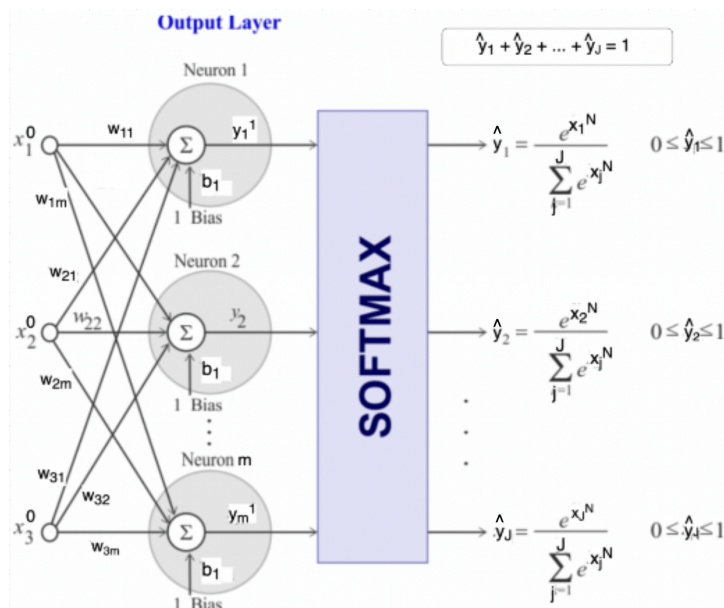


Figura 3.3: Función softmax

3.2. Entrenamiento de las redes neuronales

Originalmente partimos de un conjunto de datos, $T = \{(x_1(i), x_2(i), \dots, x_m(i), y_i)\}_{i=1, \dots, M}$, en inglés dataset, donde para cada observación hemos recogido el valor de m variables ($\{x_1, x_2, \dots, x_m\}$) a las que denominamos inputs, y la clase a la que pertenecen ($\{y\}$).

Entrenar una red neuronal es estimar el valor de los parámetros que forman parte de la red para lograr que la salida ($\{\hat{y}\}$) sea lo más parecida posible al valor real ($\{y\}$).

A la hora de entrenar, es importante dividir el dataset en dos subconjuntos de datos disjuntos. Uno de ellos es el conjunto de datos de entrenamiento, mientras que el otro es el de test. El

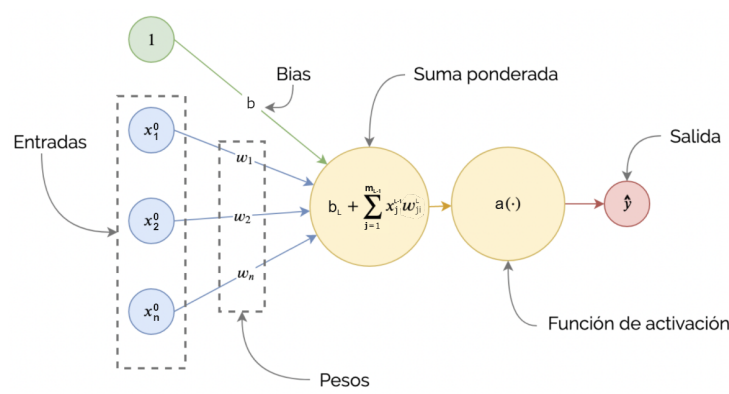


Figura 3.4: Arquitectura de una neurona

primero de ellos sirve como input de la red durante su entrenamiento, es decir, con el objetivo de estimar los pesos de las matrices W de la red de tal forma que se minimice el error de predicción. El segundo conjunto de datos, sirve para comprobar cómo de bien predice la red para unos datos nuevos no usados durante el entrenamiento. En el entrenamiento de la red debemos evaluar su rendimiento, y esta evaluación varía según el tipo de red. En nuestro caso, se trata de una red de clasificación, y por lo tanto, una forma de evaluar su rendimiento es calculando la frecuencia con la que la red clasifica correctamente los inputs, es decir, el porcentaje de número de aciertos sobre el total de elementos del conjunto de entrenamiento. A este parámetro se le denomina precisión o exactitud. Por otro lado, otra forma de evaluar su rendimiento es mediante la definición de una función de pérdida, que mide de una manera numérica la diferencia existente entre la salida ofrecida por la red y la salida correcta. Existen diversas funciones de pérdida, entre ellas el error lineal, el error cuadrático medio, entropía cruzada binaria, etc. No obstante, se va a explicar solo la que se aplica en la red neuronal creada que se presentará posteriormente.

Como se ha explicado previamente, la función de pérdida es una función que estima la desviación entre los valores reales y las predicciones hechas por la red. El objetivo de la red es minimizar esta función, es decir, reducir al mínimo la desviación entre el valor real y el valor predicho, mediante el ajuste de los parámetros de la red. En este caso, la función de pérdida utilizada es la Categorical Cross Entropy (Ec. 3.6), una función que rinde muy bien en redes de clasificación.

En cuanto a la optimización de la red, esta se lleva a cabo en dos etapas: la primera de ellas consiste en la aplicación del algoritmo de backpropagation, el cual calcula cómo afectan los valores de los parámetros al resultado de la función de pérdida, y en la segunda etapa se utiliza el algoritmo de optimización, el cual optimiza la red y cambia el valor de los parámetros a lo largo del entrenamiento para así ajustarlos. En nuestra aplicación utilizaremos como optimizador el método Adam. A continuación se explica con detalle la función de pérdida Categorical Cross

Entropy (Sec. 3.3) y el algoritmo Adam (Sec. 3.4.6).

3.3. Categorical Cross Entropy

La entropía cruzada categórica es una función de pérdida utilizada en una red neuronal para la clasificación multiclase, donde a cada muestra de los datos se le asigna una única etiqueta de clase. Esta función cuantifica la diferencia entre dos distribuciones de probabilidad discretas.

A partir de las etiquetas observadas en el conjunto de datos ($\{y_i\}_{i=1}^J$) y las predicciones obtenidas en la red ($\{\hat{y}_i\}_{i=1}^J$), con (Sec. 3.1):

$$y_i = [y_{i1}, y_{i2}, \dots, y_{iJ}] \text{ tal que}$$
$$\sum_{j=1}^J y_{ij} = 1, \quad i = 1, \dots, M, \quad y_{ij} \in \{0, 1\}$$

y

$$\hat{y}_i = [\hat{y}_{i1}, \hat{y}_{i2}, \dots, \hat{y}_{iJ}] \text{ tal que}$$
$$\sum_{j=1}^J \hat{y}_{ij} = 1, \quad i = 1, \dots, M, \quad \hat{y}_{ij} \geq 0,$$

la fórmula de la función Categorical Cross Entropy se define como:

$$C(y_i, \hat{y}_i) = - \sum_{i=1}^M \sum_{j=1}^J y_{ij} \cdot \log \hat{y}_{ij} \quad (3.6)$$

donde M es el número de inputs (tamaño de la muestra), J el número de clases, \hat{y}_{ij} es la probabilidad predicha por la red de que la muestra i -ésima pertenezca a la clase j -ésima y y_{ij} es la probabilidad real de que ocurra el evento correspondiente para la muestra i -ésima. Además, el signo menos asegura que la pérdida disminuya cuando las distribuciones se acercan más entre sí. Como $\hat{y}_{ij} \in [0, 1]$, con el fin de no calcular el logaritmo con argumento cero que no existe, realmente la biblioteca keras realiza internamente el $\log(\hat{y}_{ij} + \epsilon)$, siendo ϵ un valor muy pequeño del orden de 10^{-7} .

3.4. Adam

El algoritmo Adam [4], Adaptive Moment Estimation, es un algoritmo de optimización utilizado para estimar los pesos de la red de forma iterativa en función de los datos de entrenamiento. Su objetivo es ajustar los parámetros de la red θ (W y b) para minimizar la función de

pérdida, en nuestro caso la función Categorical Cross Entropy (Ec. 3.6). Para ello utiliza unos algoritmos de optimización, concretamente combina el algoritmo SGD con Momentum (Sec. 3.4.3) y el algoritmo RMSprop (Sec. 3.4.4) definidos en las próximos apartados. Así mismo, estos algoritmos de optimización, concretamente el Adam que es el que se va a usar en el modelo de red neuronal, utilizan el algoritmo del Backpropagation (Sec. 3.4.5) para calcular las derivadas necesarias en los algoritmos de optimización.

A continuación se va a explicar los algoritmos de optimización SGD con Momentum (Sec. 3.4.3) y RMSProp (Sec. 3.4.4) para posteriormente presentar el algoritmo Adam (Sec. 3.4.6) que combina ambos. Además, se va a explicar el algoritmo de optimización del gradiente descendente, GD, (Sec. 3.4.1) y el del gradiente descendente estocástico, SGD, (Sec. 3.4.2) ya que SGD con Momentum es una optimización de sus predecesores DG y SGD. Y para concluir, se va a presentar el algoritmo Backpropagation (Sec. 3.4.5) ya que es el usado por estos algoritmos de optimización para calcular sus derivadas.

3.4.1. Gradiente Descendiente (GD)

El método del gradiente descendente [1] es un algoritmo de optimización cuyo objetivo es converger hacia el valor mínimo de una función, $C(\theta)$, mediante un proceso iterativo.

La velocidad de convergencia del algoritmo dependerá de un parámetro α al que llamaremos tasa de aprendizaje o tamaño de paso. Una tasa demasiado grande puede impedir que se alcance la solución óptima, mientras que un valor demasiado pequeño implica que el algoritmo necesitará muchos pasos para alcanzar el óptimo. En la figura 3.5 se puede observar el efecto que tienen diferentes tasas de aprendizaje en una determinada función $C(\theta)$.

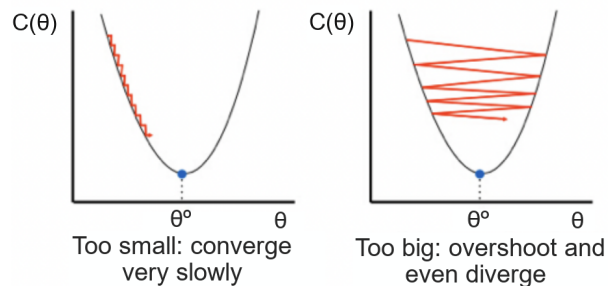


Figura 3.5: Efecto de diferentes tasas de aprendizaje

En el primer paso, el algoritmo presentado en 1 inicializa el parámetro (o vector de parámetros) θ aleatoriamente como θ_0 . A continuación se calcula el gradiente de la función respecto a θ y se evalúa en θ_0 ($g_1 \leftarrow \nabla_{\theta} C(\theta_0)$), posteriormente se actualiza el parámetro $\theta_1 \leftarrow \theta_0 - \alpha \cdot g_1$.

El algoritmo comprueba si se ha alcanzado la convergencia, es decir, si la diferencia entre el valor de los parámetros actualizados y los parámetros de la iteración anterior no superan el valor de tolerancia definido. Si es así, se ha alcanzado la convergencia y hemos llegado a un óptimo (que puede ser global o local), pero en caso contrario el algoritmo prosigue, calculando $g_t \leftarrow \nabla_{\theta} C(\theta_{t-1})$ y actualizando $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot g_t$.

Algorithm 1 Gradiente descendiente (GD)

Require: α : Tasa de aprendizaje o tamaño de paso

Require: $C(\theta)$ función objetivo con vector de parámetros θ

Require: θ_0 : vector inicial de parámetros

$t \leftarrow 0$

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} C(\theta_{t-1})$ (cálculo de los gradientes de la función de pérdida en la iteración actual t)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot g_t$ (actualización de los parámetros)

end while

return θ_t

GD en Redes Neuronales

La función a optimizar cuando trabajamos con redes neuronales es una función de pérdida, $C(\theta)$, que en nuestro caso es la Categorical Cross Entropy (Ec. 3.6), que no se evalúa solamente sobre una observación sino sobre un conjunto de observaciones (conjunto de entrenamiento) de tamaño M . En este caso, denotamos por $C_i(\theta)$ a la función de pérdida correspondiente a la i -ésima observación y la función a minimizar es:

$$C(\theta) = \sum_{i=1}^M C_i(\theta)$$

3.4.2. Gradiente Descendiente Estocástico (SGD)

El descenso de gradiente estocástico [2] es un método iterativo para optimizar una función objetivo, en nuestro caso, la función de pérdida Categorical Cross Entropy 3.6. Se trata de una optimización del algoritmo previamente definido GD 3.4.1, ya que la diferencia entre estos dos algoritmos es que en el algoritmo GD el número de muestras de entrenamiento coincide con todo el conjunto de entrenamiento, mientras que en SGD se puede elegir la cantidad de muestras a usar del conjunto de entrenamiento, las cuales son seleccionadas aleatoriamente. Aunque SGD utiliza una sola muestra como unidad de entrenamiento para entrenar más rápido, no es eficiente en conjuntos de datos muy grandes.

3.4.3. Gradiente Descendiente Estocástico con Momentum (SGD con Momentum)

Una optimización del SGD es el SGD con Momentum [7]. Este algoritmo, en vez de tomar el gradiente en el punto en particular y desplazarse en esa dirección, toma un promedio ponderado de los gradientes previos, que denotará por m_t . De esta forma, las oscilaciones en una determinada dirección se cancelan en gran medida y por lo tanto se converge mucho más rápido al valor que minimiza la función.

El algoritmo 2 inicializa aleatoriamente θ a θ_0 e inicializa la tasa de aprendizaje o tamaño de paso α y la tasa de decaimiento exponencial β , cuyos valores habituales asignados son $\alpha = 0,001$ y $\beta = 0,9$. Además, el promedio ponderado m_0 es inicializado a 0.

En cada iteración primero se calcula el gradiente mediante la derivada parcial de la función de coste C respecto de θ ($g_t \leftarrow \nabla_{\theta} C_t(\theta_{t-1})$). A continuación se actualizan los promedios ponderados $m_t \leftarrow \beta \cdot m_{t-1} + (1 - \beta) \cdot g_t$ y posteriormente se actualizan los parámetros $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot m_t$. A continuación se comprueba si el algoritmo ha alcanzado la convergencia, es decir, si la diferencia entre los parámetros actualizados y los parámetros de la iteración anterior no superan el valor de tolerancia definido. Si es así, se ha alcanzado la convergencia y hemos llegado a un óptimo (que puede ser global o local), pero en caso contrario el algoritmo prosigue, calculando $g_t \leftarrow \nabla_{\theta} C(\theta_{t-1})$ y actualizando $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot m_t$.

Algorithm 2 Gradiente Descendiente Estocástico con Momentum (SGD con Momentum)

Require: α : Tasa de aprendizaje o tamaño de paso

Require: β : Tasa de decaimiento exponencial

Require: $C(\theta)$ función objetivo con parámetros θ

Require: θ_0 : vector inicial de parámetros

$m_0 \leftarrow 0$ (Inicializar los promedios ponderados)

$t \leftarrow 0$

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} C(\theta_{t-1})$ (cálculo de los gradientes de la función de pérdida en la iteración actual t)

$m_t \leftarrow \beta \cdot m_{t-1} + (1 - \beta) \cdot g_t$ (actualización de los promedios ponderados)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot m_t$ (actualización de los parámetros)

end while

return θ_t

Además, se puede aproximar el número de elementos previos tomados en el promedio ponderado: $1/(1 - \beta)$.

3.4.4. Propagación de raíz cuadrática media (RMSProp)

RMSProp [3] también hace uso de un promedio ponderado y del gradiente de la función de pérdida respecto los parámetros, pero usa otras propiedades matemáticas, lo que permite converger más rápido al mínimo de la función. Este algoritmo trata de normalizar el gradiente con la raíz cuadrada del valor medio de los cuadrados. El objetivo es conseguir que la componente que tiene una razón de cambio mayor, avance más lento; y la que tiene razón de cambio menor avance más rápido. En este caso, se toma el promedio ponderado del cuadrado de las razones de cambio.

Algorithm 3 Propagación de raíz cuadrada media (RMSProp)

Require: α : Tasa de aprendizaje o tamaño de paso

Require: β : Tasa de decaimiento exponencial

Require: $C(\theta)$ función objetivo con parámetros θ

Require: θ_0 : vector inicial de parámetros

$v_0 \leftarrow 0$ (Inicializar los promedios ponderados)

$t \leftarrow 0$

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} C(\theta_{t-1})$ (cálculo de los gradientes de la función de pérdida en la iteración actual t)

$v_t \leftarrow \beta \cdot v_{t-1} + (1 - \beta) \cdot g_t^2$ (actualización de los promedios ponderados)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \frac{g_t}{\sqrt{v_t + \epsilon}}$ (actualización de los parámetros)

end while

return θ_t

El algoritmo 3 inicializa aleatoriamente θ a θ_0 , inicializa la tasa de aprendizaje α y la tasa de decaimiento exponencial β , cuyos valores habituales asignados son $\alpha = 0,001$ y $\beta = 0,999$. Además, el promedio ponderado m_0 es inicializado a 0. En cada iteración primero se calcula el gradiente mediante la derivada parcial de la función de coste C respecto de θ ($g_t \leftarrow \nabla_{\theta} C(\theta_{t-1})$). A continuación se actualizan los promedios ponderados $v_t \leftarrow \beta \cdot v_{t-1} + (1 - \beta) \cdot g_t^2$ y posteriormente se actualizan los parámetros $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \frac{g_t}{\sqrt{v_t + \epsilon}}$. A continuación se comprueba si el algoritmo ha alcanzado la convergencia, es decir, si la diferencia entre los parámetros actualizados y los parámetros de la iteración anterior no superan el valor de tolerancia definido. Si es así, se ha alcanzado la convergencia y hemos llegado a un óptimo (que puede ser global o local), pero en caso contrario el algoritmo prosigue, calculando $g_t \leftarrow \nabla_{\theta} C(\theta_{t-1})$ y actualizando $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \frac{g_t}{\sqrt{v_t + \epsilon}}$.

3.4.5. Backpropagation

El backpropagation [6] es un método de cálculo del gradiente utilizado en algoritmos de aprendizaje supervisado con el fin de entrenar redes neuronales. Este método usa un ciclo de propagación dividido en dos fases. En la primera de ellas, se aplica un patrón de entrada a la red y este se propaga desde la primera hasta la última a través de las capas ocultas. La salida de la red se compara con la salida esperada y se calcula un error para cada una de las salidas, y en nuestro caso el error está definido por la función de pérdida definida en la ecuación 3.6. En la segunda fase, las salidas del error se propagan hacia atrás desde la última capa hacia la capa oculta anterior. No obstante, las neuronas de la capa oculta sólo reciben una fracción de la señal total del error, basándose aproximadamente en la contribución relativa que haya aportado cada neurona a la salida original. Este proceso se repite, capa por capa, hasta que todas las neuronas de la red hayan recibido una señal de error que describa su contribución relativa al error total.

Por lo tanto, el backpropagation es el método utilizado para calcular las derivadas parciales de cada parámetro de la red con respecto al coste. Posteriormente, la optimización por la cual se va a minimizar el coste de la red se realiza mediante el algoritmo de descenso del gradiente. Este método usa el vector gradiente, un vector con las derivadas parciales de los parámetros de la red respecto al coste, y estas derivadas parciales son las calculadas con el backpropagation.

Veamos ahora cómo se calculan estas derivadas:

El objetivo es calcular cómo varía el coste (la función de pérdida) C ante un cambio del parámetro. En la red hay dos tipos de parámetros: los pesos (w) y el sesgo (b) y estos son inicializados de manera aleatoria. En el proceso de optimización necesitaremos calcular la derivada parcial respecto a w ($\frac{\partial C}{\partial w}$) y la derivada parcial respecto a b ($\frac{\partial C}{\partial b}$).

En primer lugar se calcula la derivada de los parámetros de la última capa. El superíndice indica el número de la capa a la que corresponde el parámetro. Suponiendo que N es el número total de capas de la red, recordemos que:

- $a()$ es la función de activación, siendo $a^N()$ esta función aplicada en las neuronas de la última capa (N), y $a^{N-1}()$ aplicada en la capa anterior ($N - 1$).
- $Z^N = W^N X^{N-1} + b^N = W^N a^{N-1}(Z^{N-1}) + b^N$, es resultado de la suma ponderada de la última capa, donde W^N es la matriz de pesos involucrada en el cálculo de la última capa y X^{N-1} el conjunto de las variables de estado de la capa anterior.
- $X^N = a(Z^{N-1})$
- $\hat{y}_{ij} = s(x_{ij}^N)$ (Eq. 3.5), pero para simplificar la notación en la explicación de cómo funciona el backpropagation vamos a suponer que podemos tomar $s = \text{identidad}$, es decir: $\hat{y}_{ij} = x_{ij}^N$

- $C()$ es la función de coste o de pérdida definida en la ecuación 3.6.

Como $Error = C(y, \hat{y}) = C(y, a(Z^N))$ y $Z^N = W^N a^{N-1}(Z^{N-1}) + b^N$, entonces $Error = C(a(W^N a(Z^{N-1}) + b^N))$ y aplicando la regla de la cadena obtenemos las derivadas parciales de la última capa, que se describen como:

$$\frac{\partial C}{\partial w^N} = \frac{\partial C}{\partial a^N} * \frac{\partial a^N}{\partial z^N} * \frac{\partial z^N}{\partial w^N}$$

$$\frac{\partial C}{\partial b^N} = \frac{\partial C}{\partial a^N} * \frac{\partial a^N}{\partial z^N} * \frac{\partial z^N}{\partial b^N}$$

$\frac{\partial C}{\partial a^N}$ es la derivada de la función de coste respecto de la función de activación, es decir, la derivada de la función de coste respecto al output de la red neuronal.

$\frac{\partial a^N}{\partial z^N}$ es la derivada de la función de activación respecto a la suma ponderada de la neurona, es decir, la derivada de la función de activación.

$\frac{\partial z^N}{\partial w^N}$ y $\frac{\partial z^N}{\partial b^N}$ son las derivadas de la suma ponderada al variar los pesos o el sesgo respectivamente. Por lo tanto, $\frac{\partial z^N}{\partial b^N} = 1$ ya que b es una variable independiente, y $\frac{\partial z^N}{\partial w^N} = a_i^{N-1}$ es el valor de entrada a la neurona.

Además, $\frac{\partial C}{\partial z^N} = \frac{\partial C}{\partial a^N} * \frac{\partial a^N}{\partial z^N}$ representa cómo varía el error en función de la suma ponderada, es decir, en qué grado se modifica el error al realizar un cambio en la suma de la neurona. Si la derivada es grande, significa que ante un pequeño cambio en el valor de la neurona, este se verá reflejado en el resultado final. No obstante, si la derivada es pequeña significa que el valor de la suma no afectará al valor de la red. Entonces, esta derivada indica qué responsabilidad tiene la neurona en el resultado final y por tanto en el error. Además, recibe el nombre de error imputado a la neurona, y se representa como ζ^N .

Por lo tanto, $\frac{\partial C}{\partial w^N} = \zeta^N a_i^{N-1}$ y $\frac{\partial C}{\partial b^N} = \zeta^N$.

De esta forma, hemos calculado el error para cada una de las neuronas de la última capa (ζ^N) y el error para cada una de las derivadas parciales ($\frac{\partial C}{\partial w^N}$ y $\frac{\partial C}{\partial b^N}$). No obstante, todavía se necesita seguir derivando para poder calcular el resto de derivadas de la red. Para ello, realizamos el mismo cálculo para la capa anterior, la capa $N - 1$.

Como $Error = C(a^N(Z^N(Z^{N-1})))$ y $Z^{N-1} = W^{N-1} a^{N-2}(Z^{N-2}) + b^{N-1}$, entonces $Error = C(a^N(Z^N(W^{N-1} a^{N-2}(Z^{N-2}) + b^{N-1}) + b^N))$, y aplicando la regla de la cadena obtenemos las derivadas parciales de la penúltima capa se describen como:

$$\frac{\partial C}{\partial w^{N-1}} = \frac{\partial C}{\partial a^N} * \frac{\partial a^N}{\partial z^N} * \frac{\partial z^N}{\partial a^{N-1}} * \frac{\partial a^{N-1}}{\partial z^{N-1}} * \frac{\partial z^{N-1}}{\partial w^{N-1}}$$

$$\frac{\partial C}{\partial b^{N-1}} = \frac{\partial C}{\partial a^N} * \frac{\partial a^N}{\partial z^N} * \frac{\partial z^N}{\partial a^{N-1}} * \frac{\partial a^{N-1}}{\partial z^{N-1}} * \frac{\partial z^{N-1}}{\partial b^{N-1}}$$

donde $\frac{\partial C}{\partial a^N} \frac{\partial a^N}{\partial z^N} = \zeta^N$ y por lo tanto $\frac{\partial C}{\partial a^N} * \frac{\partial a^N}{\partial z^N} * \frac{\partial z^N}{\partial a^{N-1}} * \frac{\partial a^{N-1}}{\partial z^{N-1}} = \frac{\partial C}{\partial z^{N-1}} = \zeta^{N-1}$. Además, $\frac{\partial z^{N-1}}{\partial w^{N-1}} = a^{N-2}$ y $\frac{\partial z^{N-1}}{\partial b^{N-1}} = 1$.

De esta forma, tenemos calculadas:

- El error de la última capa: $\zeta^N = \frac{\partial C}{\partial a^N} * \frac{\partial a^N}{\partial z^N}$
- El error de la capa anterior: $\zeta^{N-1} = W^N \zeta^N * \frac{\partial a^{N-1}}{\partial z^{N-1}}$
- Las derivadas de la capa usando el error: $\frac{\partial C}{\partial b^{N-1}} = \zeta^{N-1}$ y $\frac{\partial C}{\partial w^{N-1}} = \zeta^{N-1} * a^{N-2}$

Y así sucesivamente recorriendo todas las capas de la red. Por lo tanto, con estas cuatro expresiones, con un único pase se pueden calcular todos los errores y derivadas parciales de la red.

3.4.6. Algoritmo Adam

En el algoritmo Adam (4) se mantiene una tasa de aprendizaje para cada parámetro de la red, es decir, el tamaño de paso del gradiente descendiente estocástico se ajusta a medida que se desarrolla el aprendizaje.

Los parámetros de configuración de Adam son:

- α : se conoce como tasa de aprendizaje o tamaño de paso. La proporción en la que se actualizan las ponderaciones. Los valores más grandes dan como resultado un aprendizaje inicial más rápido antes de que se actualice la tasa. Los valores más pequeños ralentizan el aprendizaje durante el entrenamiento
- β_1 : tasa de caída exponencial para las estimaciones del primer momento. Suele ser de 0,9.
- β_2 : tasa de caída exponencial para las estimaciones del segundo momento. Suele ser de 0,999.
- ϵ : número muy pequeño para evitar cualquier división por cero en la implementación (por ejemplo, 10^{-9}).

Como hemos dicho anteriormente, Adam combina SGD Momentum con RMSProp, y su algoritmo se define como:

Algorithm 4 Adam

Require: α : Tasa de aprendizaje o tamaño de paso

Require: $\beta_1, \beta_2 \in [0, 1)$: Tasas de decaimiento exponenciales para los momentos estimados

Require: $C(\theta)$ función objetivo con parámetros θ

Require: θ_0 : vector inicial de parámetros

$m_0 \leftarrow 0$ (Inicializar el primer momento)

$v_0 \leftarrow 0$ (Inicializar el segundo momento)

$t \leftarrow 0$

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} C_t(\theta_{t-1})$ (cálculo de los gradientes de la función de pérdida en la iteración actual t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (actualización de los primeros momentos estimados)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (actualización de los segundos momentos estimados)

$\hat{m}_t \leftarrow \frac{m_t}{1 - \beta_1^t}$ (cálculo de la corrección de los primeros momentos)

$\hat{v}_t \leftarrow \frac{v_t}{1 - \beta_2^t}$ (cálculo de la corrección de los segundos momentos)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$ (actualización de los parámetros)

end while

return θ_t

Como podemos observar el algoritmo Adam combina SGD con Momentum y RMSProp. Como se ha explicado anteriormente, m_t y v_t son los promedios ponderados que se calculaban en los algoritmos SGD con Momentum y RMSProp respectivamente. m_t y v_t representan los dos momentos, m_t es el que modela la media de los gradientes a lo largo del tiempo y v_t hace lo mismo pero modelando la varianza; y lo mismo ocurre respectivamente con el sesgo. En este algoritmo, en las primeras iteraciones tiene un arranque lento, es decir, en las actualizaciones de los parámetros de las primeras iteraciones se tiene muy poco en cuenta el valor de los promedios ponderados. Para solventar este problema en cada iteración se realiza una corrección de los momentos antes de actualizar los parámetros.

Concretamente, el algoritmo calcula un promedio móvil exponencial del gradiente y el gradiente al cuadrado, y los parámetros β_1 y β_2 controlan las tasas de caída de estos promedios móviles.

El algoritmo 4 se inicia tomando unos valores aleatorios de los parámetros θ (W y b), que al ser los iniciales los denominamos como θ_0 . Además, inicializa la tasa de aprendizaje o tamaño de paso α y las tasas de decaimiento exponenciales β_1 y β_2 . Una buena configuración determinada para este algoritmo es definir α como 0.001, β_1 como 0.9 y β_2 como 0.999. Además, los primeros y segundos momentos m_0 y v_0 son inicializados a 0. En cada iteración primero se calcula el gradiente mediante la derivada parcial de la función de coste C respecto de los parámetros θ ($g_t \leftarrow \nabla_{\theta} C_t(\theta_{t-1})$) resultantes de la iteración anterior. A continuación se actualizan los primeros y segundos momentos mediante las expresiones $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ y

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ y posteriormente se calcula la corrección de los primeros y segundos momentos mediante las expresiones $\hat{m}_t \leftarrow \frac{m_t}{1 - \beta_1^t}$ y $\hat{v}_t \leftarrow \frac{v_t}{1 - \beta_2^t}$. Después, se actualizan los parámetros de esta iteración mediante la expresión $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$. A continuación se comprueba si el algoritmo ha alcanzado la convergencia, es decir, si la diferencia entre los parámetros actualizados y los parámetros de la iteración anterior superan el valor de tolerancia definido. Si es así, el algoritmo prosigue, pero en caso contrario esto implica que ha alcanzado la convergencia y termina.

3.4.7. Ejemplo

Una vez visto como funciona Adam, veamos sobre un caso particular cómo este algoritmo calcula las derivadas de los parámetros de una red para ajustarlos con el fin de minimizar una función de pérdida. Concretamente, se han realizado los cálculos de estas derivadas para una sola observación, suponiendo que la función a minimizar es $C(\theta) = C_i(\theta)$, pero el desarrollo sería análogo en el caso de $C(\theta) = \frac{1}{M} \sum_{i=1}^N C_i(\theta)$. En este ejemplo, la arquitectura de la red consiste en una capa de entrada, una capa oculta densa y una capa de salida, cada una de ellas con dos neuronas. Además, la función de activación usada en la capa oculta y en la capa de salida es distinta; en la primera de ellas se usa la función ReLU (ec. 3.2) y en la segunda de ellas se usa la función softmax (ec. 3.5). Podemos observar esta red en la imagen 3.6.

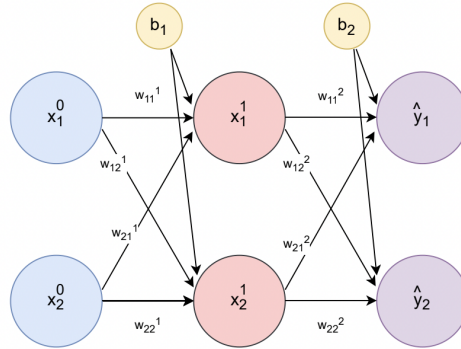


Figura 3.6: Arquitectura de la red del ejemplo práctico

- $x_1^1 = a(z_1^1) = a(x_1^0 \cdot w_{11}^1 + x_2^0 \cdot w_{21}^1 + b_1) = \max(0, x_1^0 \cdot w_{11}^1 + x_2^0 \cdot w_{21}^1 + b_1)$
- $x_2^1 = a(z_2^1) = a(x_1^0 \cdot w_{12}^1 + x_2^0 \cdot w_{22}^1 + b_1) = \max(0, x_1^0 \cdot w_{12}^1 + x_2^0 \cdot w_{22}^1 + b_1)$
- $\hat{y}_1 = s(z_1^2) = s(x_1^1 \cdot w_{11}^2 + x_2^1 \cdot w_{21}^2 + b_2)$
- $\hat{y}_2 = s(z_2^2) = \frac{e^{z_1^2}}{e^{z_1^2} + e^{z_2^2}}$

- $\hat{y}_2 = s(z_2^2) = s(x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2)$
- $\hat{y}_2 = s(z_2^2) = \frac{e^{z_2^2}}{e^{z_1^2} + e^{z_2^2}}$
- $C(y_i, \hat{y}_i) = - \sum_{i=1}^J y_i \cdot \log \hat{y}_i = y_1 \cdot \log \hat{y}_1 + y_2 \cdot \log \hat{y}_2$

- CÁLCULO DE LAS DERIVADAS PARCIALES DE LA FUNCIÓN DE COSTE C RESPECTO A LOS PARÁMETROS DE LA ÚLTIMA CAPA (b_2 , w_{11}^2 , w_{12}^2 , w_{21}^2 y w_{22}^2):

Aplicamos las siguientes fórmulas ya definidas en la sección 3.4.5: $\frac{\partial C}{\partial w^N} = \frac{\partial C}{\partial a^N} \cdot \frac{\partial a^N}{\partial z^N} \cdot \frac{\partial z^N}{\partial w^N}$
y $\frac{\partial C}{\partial b^N} = \frac{\partial C}{\partial a^N} \cdot \frac{\partial a^N}{\partial z^N} \cdot \frac{\partial z^N}{\partial b^N}$

- $\frac{\partial z^N}{\partial w^N}$
 - $\frac{\partial z_1^2}{\partial w_{11}^2} = x_1^2$
 - $\frac{\partial z_1^2}{\partial w_{21}^2} = x_1^2$
 - $\frac{\partial z_2^2}{\partial w_{12}^2} = x_1^2$
 - $\frac{\partial z_2^2}{\partial w_{22}^2} = x_1^2$
- $\frac{\partial z^N}{\partial b^N}$
 - $\frac{\partial z_1^2}{\partial b^2} = 1$
 - $\frac{\partial z_2^2}{\partial b^2} = 1$
- $\frac{\partial a^N}{\partial z^N}$
 - $\frac{\partial \hat{y}_1}{\partial z_1^2} = \frac{e^{z_1^2} \cdot (2 \cdot e^{z_1^2} + e^{z_2^2})}{(e^{z_1^2} + e^{z_2^2})^2}$
 - $\frac{\partial \hat{y}_2}{\partial z_2^2} = \frac{e^{z_2^2} \cdot (2 \cdot e^{z_2^2} + e^{z_1^2})}{(e^{z_1^2} + e^{z_2^2})^2}$
- $\frac{\partial C}{\partial a^N}$
 - $\frac{\partial C}{\partial \hat{y}_1} = \frac{y_1}{\hat{y}_1}$
 - $\frac{\partial C}{\partial \hat{y}_2} = \frac{y_2}{\hat{y}_2}$

Por lo tanto, juntando las derivadas parciales, sustituyendo y expresando en función de x y w obtenemos:

- $\frac{\partial C}{\partial w_{11}^2} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial w_{11}^2} = y_1 \cdot x_1^2 \cdot \frac{2 \cdot e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}{e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}$
- $\frac{\partial C}{\partial w_{12}^2} = \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial w_{12}^2} = y_2 \cdot x_1^2 \cdot \frac{2 \cdot e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2} + e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2}}{e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}$
- $\frac{\partial C}{\partial w_{21}^2} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial w_{21}^2} = y_1 \cdot x_2^2 \cdot \frac{2 \cdot e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}{e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}$
- $\frac{\partial C}{\partial w_{22}^2} = \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial w_{22}^2} = y_2 \cdot x_2^2 \cdot \frac{2 \cdot e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2} + e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2}}{e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}$
- $\frac{\partial C}{\partial b_2} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial b_2} + \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial b_2} = y_1 \cdot \frac{2 \cdot e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}{e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}} + y_2 \cdot \frac{2 \cdot e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2} + e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2}}{e^{x_1^2 \cdot w_{11}^2 + x_2^2 \cdot w_{21}^2 + b_2} + e^{x_1^2 \cdot w_{12}^2 + x_2^2 \cdot w_{22}^2 + b_2}}$

- CÁLCULO DE LAS DERIVADAS PARCIALES DE LA FUNCIÓN DE COSTE C RESPECTO A LOS PARÁMETROS DE LA CAPA OCULTA (b_2 , w_{11}^2 , w_{12}^2 , w_{21}^2 y w_{22}^2):

Aplicamos las siguientes fórmulas ya definidas en la sección 3.4.5:

$$\frac{\partial C}{\partial w^{N-1}} = \frac{\partial C}{\partial a^N} \cdot \frac{\partial a^N}{\partial z^N} \cdot \frac{\partial z^N}{\partial a^{N-1}} \cdot \frac{\partial a^{N-1}}{\partial z^{N-1}} \cdot \frac{\partial z^{N-1}}{\partial w^{N-1}}$$

y

$$\frac{\partial C}{\partial b^{N-1}} = \frac{\partial C}{\partial a^N} \cdot \frac{\partial a^N}{\partial z^N} \cdot \frac{\partial z^N}{\partial a^{N-1}} * \frac{\partial a^{N-1}}{\partial z^{N-1}} \cdot \frac{\partial z^{N-1}}{\partial b^{N-1}}$$

- $\frac{\partial z^{N-1}}{\partial w^{N-1}}$
 - $\frac{\partial z_1^1}{\partial w_{11}^1} = x_1^0$
 - $\frac{\partial z_1^1}{\partial w_{21}^1} = x_2^0$
 - $\frac{\partial z_2^1}{\partial w_{12}^1} = x_1^0$
 - $\frac{\partial z_2^1}{\partial w_{22}^1} = x_2^0$
- $\frac{\partial a^{N-1}}{\partial z^{N-1}}$
 - $\frac{\partial x_1^1}{\partial z_1^1} = 0$ si $z_1^1 = 0$, $\frac{\partial x_1^1}{\partial z_1^1} = 1$ si $z_1^1 \neq 0$
 - $\frac{\partial x_2^1}{\partial z_2^1} = 0$ si $z_2^1 = 0$, $\frac{\partial x_2^1}{\partial z_2^1} = 1$ si $z_2^1 \neq 0$
- $\frac{\partial z^N}{\partial a^{N-1}}$
 - $\frac{\partial z_1^2}{\partial x_1^1} = w_{11}^2$

- $\frac{\partial z_1^2}{\partial x_2^2} = w_{21}^2$
- $\frac{\partial z_2^2}{\partial x_1^2} = w_{12}^2$
- $\frac{\partial z_2^2}{\partial x_2^2} = w_{22}^2$
- $\frac{\partial a^N}{\partial z^N}$
 - $\frac{\partial \hat{y}_1}{\partial z_1^2} = \frac{e^{z_1^2} \cdot (2 \cdot e^{z_1^2} + e^{z_2^2})}{(e^{z_1^2} + e^{z_2^2})^2}$
 - $\frac{\partial \hat{y}_2}{\partial z_2^2} = \frac{e^{z_2^2} \cdot (2 \cdot e^{z_2^2} + e^{z_1^2})}{(e^{z_1^2} + e^{z_2^2})^2}$
- $\frac{\partial C}{\partial a^N}$
 - $\frac{\partial C}{\partial \hat{y}_1} = \frac{y_1}{\hat{y}_1}$
 - $\frac{\partial C}{\partial \hat{y}_2} = \frac{y_2}{\hat{y}_2}$

Para hacer los siguientes cálculos suponemos que $z_1^1 \neq 0$ y $z_2^1 \neq 0$, ya que en caso contrario las derivadas serían 0. Por lo tanto, juntando las derivadas parciales, sustituyendo y expresando en función de x y w obtenemos:

- $\frac{\partial C}{\partial w_{11}^1} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial x_1^2} \cdot \frac{\partial x_1^1}{\partial z_1^1} \cdot \frac{\partial z_1^1}{\partial w_{11}^1} + \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial x_1^2} \cdot \frac{\partial x_1^1}{\partial z_1^1} \cdot \frac{\partial z_1^1}{\partial w_{11}^1} = x_1^0 \cdot [y_1 \cdot w_{11}^2 \cdot \frac{2 \cdot e^{z_1^2} + e^{z_2^2}}{(e^{z_1^2} + e^{z_2^2})^2} + y_2 \cdot w_{12}^2 \cdot \frac{2 \cdot e^{z_2^2} + e^{z_1^2}}{(e^{z_1^2} + e^{z_2^2})^2}]$
- $\frac{\partial C}{\partial w_{12}^1} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial x_2^2} \cdot \frac{\partial x_2^1}{\partial z_2^1} \cdot \frac{\partial z_2^1}{\partial w_{12}^1} + \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial x_2^2} \cdot \frac{\partial x_2^1}{\partial z_2^1} \cdot \frac{\partial z_2^1}{\partial w_{12}^1} = x_1^0 \cdot [y_1 \cdot w_{21}^2 \cdot \frac{2 \cdot e^{z_1^2} + e^{z_2^2}}{(e^{z_1^2} + e^{z_2^2})^2} + y_2 \cdot w_{22}^2 \cdot \frac{2 \cdot e^{z_2^2} + e^{z_1^2}}{(e^{z_1^2} + e^{z_2^2})^2}]$
- $\frac{\partial C}{\partial w_{21}^1} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial x_1^2} \cdot \frac{\partial x_1^1}{\partial z_1^1} \cdot \frac{\partial z_1^1}{\partial w_{21}^1} + \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial x_1^2} \cdot \frac{\partial x_1^1}{\partial z_1^1} \cdot \frac{\partial z_1^1}{\partial w_{21}^1} = x_2^0 \cdot [y_1 \cdot w_{11}^2 \cdot \frac{2 \cdot e^{z_1^2} + e^{z_2^2}}{(e^{z_1^2} + e^{z_2^2})^2} + y_2 \cdot w_{12}^2 \cdot \frac{2 \cdot e^{z_2^2} + e^{z_1^2}}{(e^{z_1^2} + e^{z_2^2})^2}]$
- $\frac{\partial C}{\partial w_{22}^1} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial x_2^2} \cdot \frac{\partial x_2^1}{\partial z_2^1} \cdot \frac{\partial z_2^1}{\partial w_{22}^1} + \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial x_2^2} \cdot \frac{\partial x_2^1}{\partial z_2^1} \cdot \frac{\partial z_2^1}{\partial w_{22}^1} = x_2^0 \cdot [y_1 \cdot w_{21}^2 \cdot \frac{2 \cdot e^{z_1^2} + e^{z_2^2}}{(e^{z_1^2} + e^{z_2^2})^2} + y_2 \cdot w_{22}^2 \cdot \frac{2 \cdot e^{z_2^2} + e^{z_1^2}}{(e^{z_1^2} + e^{z_2^2})^2}]$
- $\frac{\partial C}{\partial b_1} = \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial x_1^1} \cdot \frac{\partial x_1^1}{\partial z_1^1} + \frac{\partial C}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_1^2} \cdot \frac{\partial z_1^2}{\partial x_2^1} \cdot \frac{\partial x_2^1}{\partial z_2^1} + \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} + \frac{\partial C}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial z_2^2} \cdot \frac{\partial z_2^2}{\partial x_2^1} \cdot \frac{\partial x_2^1}{\partial z_2^1} = y_1 \cdot w_{11}^2 \cdot \frac{2 \cdot e^{z_1^2} + e^{z_2^2}}{(e^{z_1^2} + e^{z_2^2})^2} + y_1 \cdot w_{11}^2 \cdot \frac{2 \cdot e^{z_1^2} + e^{z_2^2}}{(e^{z_1^2} + e^{z_2^2})^2} + y_2 \cdot w_{12}^2 \cdot \frac{2 \cdot e^{z_2^2} + e^{z_1^2}}{(e^{z_1^2} + e^{z_2^2})^2} + y_2 \cdot w_{22}^2 \cdot \frac{2 \cdot e^{z_2^2} + e^{z_1^2}}{(e^{z_1^2} + e^{z_2^2})^2}$

Capítulo 4

Resultados

Una vez estudiados los fundamentos teóricos de las redes neuronales más sencillas, vamos a aplicarlos de forma práctica en un desafío propuesto por la empresa donde he desarrollado mis prácticas curriculares.

Esta empresa es IronArray y su actividad consiste en proporcionar herramientas para llevar a cabo una compresión y descompresión de datos rápida y eficaz, adaptándose a las necesidades de rendimiento y almacenamiento del cliente. Una de las principales bibliotecas desarrolladas por el equipo es C-Blosc2 [8], un compresor y metacompresor al mismo tiempo. Entre los distintos métodos de compresión usados por la empresa, encontramos tres parámetros a ajustar.

El primero de ellos es el *filtro*, un programa que reordena los datos sin modificar su tamaño, de manera que el tamaño inicial y final sean iguales. Un filtro consta de codificador y decodificador. Partimos de que un dato es un conjunto o bloque de bytes de cierto tamaño. Por lo tanto, los filtros reorganizan los bytes o bits del bloque de datos a comprimir, acumulando a la izquierda los bytes o bits que son más significativos de cada dato y los bytes o bits menos significativos a la derecha. En una representación binaria (ceros y unos) donde cada número representa un byte, los bytes más significativos son los que tienen el mayor valor y están situados más a la izquierda, mientras que los menos significativos son los que tienen el menor valor y se sitúan más a la derecha. De entre los filtros que proporciona el compresor C-Blosc2, vamos a distinguir entre shuffle, el cual reordena los datos a nivel de byte; y bitshuffle, que los reordena a nivel de bit. En la figura 4.1 se puede apreciar el esquema del filtro shuffle, siendo el del bitshuffle el mismo pero en vez de tratar con bytes lo hace con bits. El segundo parámetro es el códec, un algoritmo capaz de comprimir y descomprimir un flujo de datos con el objetivo de reducir su tamaño y permitir una transmisión de datos más rápida. Disponemos de cinco códecs distintos: blosclz, zlib, zstd, lz4 i lz4hc. El codificador de filtro se aplica antes del compresor de códec (o codificador de códec) para que los datos sean más fáciles de comprimir y el decodificador

de filtro se usa después del descompresor de códec (o decodificador de códec) para restaurar la disposición original de los datos. El tercero de ellos es el splitmode, el cual especifica si se va a dividir los datos en splits o bloques (conjuntos de splits). Si el splitmode es split, significa que divide el dataset en splits, pero si el splitmode es nosplit, entonces se divide en bloques.

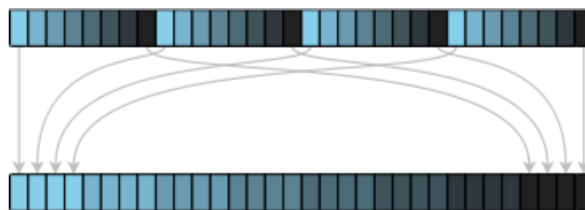


Figura 4.1: Esquema del filtro shuffle, en la parte superior aparecen 4 datos de 8 bytes, mientras que en la parte inferior aparecen los bytes filtrados y reordenados.

Una de las variables a definir es el favour, es decir, el parámetro en el que el usuario decide si quiere maximizar el ratio de compresión (max_cratio), la velocidad de compresión (max_speed) o equilibrio entre los dos (balance). Dependiendo del favour seleccionado, se le aplicará un determinado peso a la función de puntuación de categorías definida previamente. Por lo tanto, para cada favour se ha creado un modelo de red neuronal concreto.

Por experiencia previa, sabemos que cuando se quiere maximizar el ratio de compresión, los códecs más efectivos son zlib, zstd y lz4hc; cuando lo que se quiere maximizar es la velocidad de compresión, son idóneos los códecs blosclz y lz4; y finalmente para optimizar el modo balance, los mejores resultados los proporcionan los códecs blosclz y zstd, el primero de ellos en modo split, y el segundo en modo nosplit. Por lo tanto, dada esta información se ha decidido fijar uno de los filtros disponibles, concretamente el shuffle, y lo mismo en los dos primeros tipos de balance con unos de los splitmode disponibles, el split. En las siguientes tablas se muestra las combinaciones de parámetros (categorías) entre las cuales cada red neuronal deberá clasificar.

Favour
max_cratio
max_speed
balance

Cuadro 4.1: Tipos de favour según objetivo a optimizar

Además del compresor con estos parámetros, ironArray dispone de un programa llamado entropy_probe, el cual dado un dataset calcula su entropía, es decir, el grado de desorden de los datos de ese dataset. Con el objetivo de obtener información más detallada sobre los datos comprimidos, he realizado unas implementaciones sobre él. Como consecuencia, al programa se le pasa un dataset y su función es hacer una estimación del ratio y velocidad de compresión para un determinado filtro, aunque en los modelos de red neuronal creados se ha decidido

Filtro	Códec	Splitmode	Identificador de categoría	Categoría
shuffle	zlib	split	0	zlib-shuffle-split
shuffle	zstd	split	1	zstd-shuffle-split
shuffle	lz4hc	split	2	lz4hc-shuffle-split

Cuadro 4.2: Categorías de la red neuronal para optimizar el cratio (max_cratio)

Filtro	Códec	Splitmode	Identificador de categoría	Categoría
shuffle	blosclz	split	0	blosclz-shuffle-split
shuffle	lz4	split	1	lz4-shuffle-split

Cuadro 4.3: Categorías de la red neuronal para optimizar el speed (max_speed)

fijar como filtro el shuffle. Este estimador funciona siempre con el splitmode en modo split. Posteriormente, estos ratios y velocidades serán los inputs del modelo de red neuronal de multi clasificación creado.

El objetivo es optimizar estos parámetros de compresión para cada dataset que se vaya a comprimir. Para ello, se ha creado un modelo de red neuronal que dados el ratio y velocidad de compresión de unos datos, predice cuál es la mejor combinación que se le puede aplicar. Concretamente, se ha creado un modelo de red neuronal para cada objetivo posible de favour; uno para optimizar el ratio de compresión, otro para la velocidad de compresión y otro para optimizar el balance. Los inputs de las redes son los ratios y velocidades de compresión de entropía, y los outputs son las combinaciones óptimas de parámetros, a las que denominamos categorías. Por lo tanto, estos modelos consisten en redes de multi clasificación.

4.1. Obtención de los datos

Para entrenar la red neuronal necesitamos unos inputs. En este caso, estos son el ratio y la velocidad de compresión. El ratio de compresión es la relación entre el tamaño del dato comprimido y el tamaño del dato sin comprimir. En cuanto a la velocidad de compresión, esta se refiere al tiempo (Bytes/segundo) que se tarda en comprimir los datos.

Como se ha introducido anteriormente, con el fin de obtener estos inputs de los datos he implementado unas funciones en el programa en lenguaje C denominado entropy_probe.c, el cual ha sido proporcionado por ironArray. Como consecuencia, al ejecutar el programa para un determinado dataset, se obtiene un documento csv donde para cada porción del dataset (split), obtenemos información sobre su ratio de compresión y velocidad de compresión, entre otros. Esta información se obtiene por una parte para los filtros y codecs seleccionados, y por otra parte, si se especifica en las configuraciones de ejecución del programa, para un códec especial

Filtro	Códec	Splitmode	Identificador de categoría	Categoría
shuffle	blosclz	split	0	blosclz-shuffle-split
shuffle	lz4	nosplit	1	lz4-shuffle-nosplit

Cuadro 4.4: Categorías de la red neuronal para optimizar el balance

desarrollado por ironArray. Este códec se denomina `entropy_mode`. Su finalidad es medir la entropía del dataset, es decir, el grado de ordenación de los datos dentro del dataset. Además, en este caso sirve como estimador de ratio y velocidad de compresión en la red neuronal. En los modelos de red neuronal creados, se ha utilizado tres datasets. El primero de ellos contiene información sobre precipitaciones, el segundo sobre temperaturas y el tercero sobre el viento. Para la creación de las redes neuronales se ha utilizado un dataset resultado de la unión de los tres nombrados.

A continuación se muestra las gráficas donde aparecen las nubes de puntos para cada clase, es decir, el ratio y velocidad de compresión para cada muestra de los tres datasets.

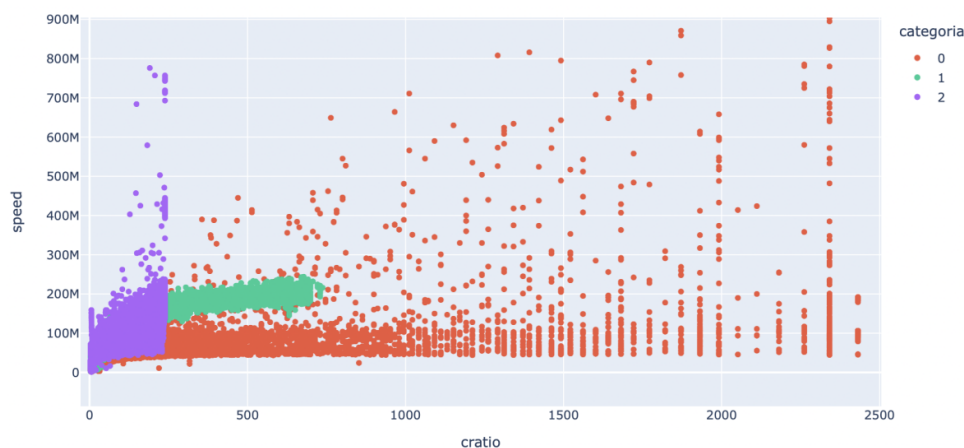


Figura 4.2: Mediciones del dataset para el favour en modo `max_cratio`.

En la figura 4.2 podemos observar las mediciones de ratio y velocidad de compresión para los datos del dataset global, y lo mismo en las figuras 4.3 y 4.4; diferenciándose entre ellas por el modo de favour seleccionado.

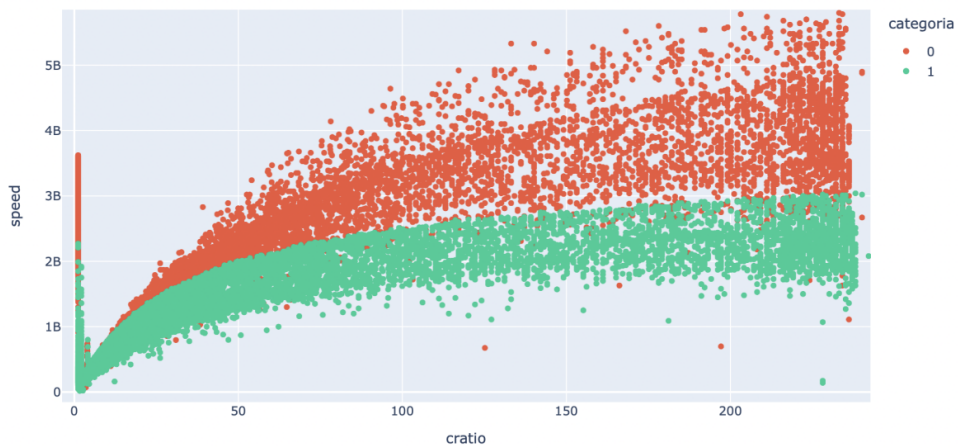


Figura 4.3: Mediciones del dataset para el favour en modo max_speed.

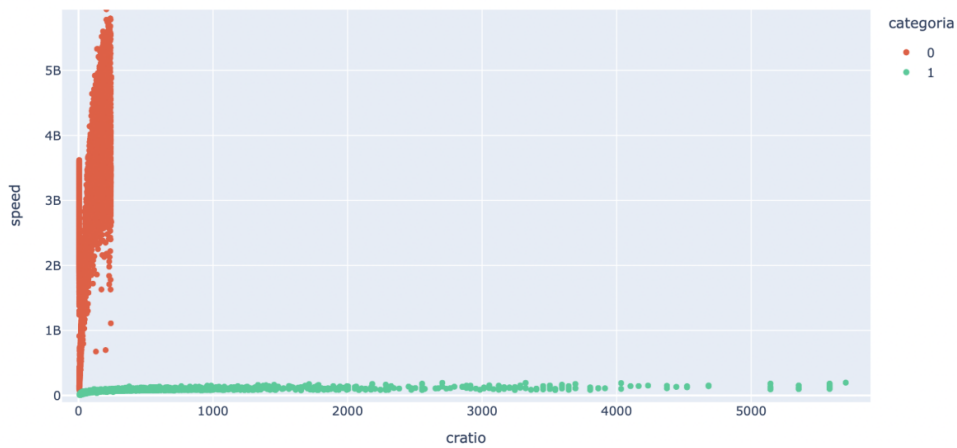


Figura 4.4: Mediciones del dataset para el favour en modo balance.

4.2. Preparación y manipulación de los datos

Una vez obtenidos los archivos csv con la información sobre los datasets, la siguiente tarea a realizar es la preparación de estos datos para convertirlos en los inputs definitivos de los modelos de red neuronal.

Esta preparación se ha llevado a cabo en el Jupyter Notebook, una aplicación web libre para crear y compartir documentos que contienen bloques de código con sus correspondientes salidas

por pantalla, y al mismo tiempo bloques de texto potenciados con el lenguaje Markdown, de tal forma que soporta LaTeX y HTML. Es muy útil en el análisis de datos ya que permite llevar un registro de las técnicas que se han empleado, además de que permite generar informes en HTML y pdf. Además, para realizar ciertas operaciones de manejo de datos se han utilizado las librerías pandas, matplotlib, pyplot, plotly y numpy.

El objetivo de la red neuronal es que, dado un dataset, clasifique cual es la combinación óptima de filtro y códec para cada split. De esta manera, una vez se tienen las predicciones para todos los splits que forman el dataset, se aplicaría a todo el dataset la combinación de parámetros de compresión que la red haya predicho con mayor frecuencia.

Como he descrito anteriormente, los inputs de la red son el ratio y velocidad de compresión de entropía. A cada par de estos inputs de un determinado split se le debe asignar una etiqueta, que en las redes neuronales de multi clasificación se refiere a la categoría óptima a la que pertenece ese par de inputs.

Por otra parte, en los datasets puede haber un tipo especial de valor, los denominados special values. Estos valores son conjuntos de bits que, al ser reordenados por alguno de los dos filtros, se convierten en agrupaciones de un mismo número. Como consecuencia, su compresión es prácticamente inmediata y por lo tanto producen unos valores de ratio y velocidad de compresión extremadamente elevados. Por consiguiente, se ha decidido eliminar estos valores del conjunto de datos a utilizar como input de la red neuronal.

Una vez realizado este filtrado, se han agrupado los datos según el filtro usado. A continuación se ha procedido a calcular cuál es la categoría que produce mejores resultados para cada split. Para ello se ha definido una función de puntuación (score), la cual se trata de maximizar. Esta función consiste en sumar el ratio y velocidad de compresión normalizados:

$$S(\text{cratio}_i, \text{speed}_i) = \text{theta} * \frac{\text{cratio}_i - \text{mean_cratios}}{\text{std_cratios}} + (1 - \text{theta}) * \frac{\text{speed}_i - \text{mean_speeds}}{\text{std_speeds}} \quad (4.1)$$

para $i = 1, 2, \dots, N$ siendo N el número de splits del dataset. El peso theta dependerá del favour seleccionado.

Una vez calculado el score de las categorías de cada filtro, seleccionamos la mayor de ellas ya que a mayor ratio y velocidad, mejor es la compresión. Esta categoría seleccionada como la mejor pasa a ser la etiqueta que se le introducirá a la red, la cual estará relacionada con el ratio y velocidad de compresión del entropy_mode.

4.3. Construcción del modelo de red neuronal

Una vez preparados los inputs de la red, el siguiente paso es crear y construir la red. Para ello he utilizado las librerías keras y tensorflow.

4.3.1. Arquitectura de la red

Las redes neuronales creadas están formadas por una capa de entrada, una capa oculta y otra capa de salida. La capa de entrada dispone de dos neuronas, una para cada input; la capa oculta tiene 2 neuronas, y finalmente la capa de salida tiene cuatro neuronas, ya que la clasificación se realiza entre cuatro categorías. En cuanto a las funciones de optimización, en todas sus capas menos la de salida se ha usado la función ReLu (3.2). Además, en la capa de salida la función de activación utilizada es la función softmax (3.5).

4.3.2. Compilación de la red

Una vez creada la arquitectura de las redes, el siguiente paso es compilarlas. Para ello hay que definir tres parámetros: optimizador, función de pérdida y métricas. El optimizador es el mecanismo por el cual el modelo se actualiza basándose en los datos de entrenamiento que ve, con el objetivo de mejorar su rendimiento. La función de pérdida expresa cómo el modelo será capaz de medir su rendimiento en los datos de entrenamiento, y por lo tanto cómo será capaz de dirigirse en la dirección correcta.

En este caso, el optimizador utilizado es Adam, el cual ha sido descrito anteriormente. La función de pérdida empleada es la Categorical Cross Entropy, que también ha sido explicada en el capítulo 3.3. Finalmente, en cuanto a las métricas a monitorizar durante el entrenamiento, existen diversas pero la escogida es la precisión (accuracy), es decir, la fracción de datos que han sido clasificados correctamente.

4.3.3. Adecuación del modelo (fitting)

Tras haber construido y compilado las redes, el último paso es la adecuación. En este paso introducimos los inputs y etiquetas de entrenamiento y de test, especificando también las épocas que queremos que se realicen en la red. Ahora ya se puede proceder al entrenamiento de las redes neuronales.

4.3.4. Resultados de la red

Una vez entrenadas, las redes neuronales podemos observar su rendimiento y su función de pérdida representadas mediante gráficas. Como ya se ha explicado, la precisión es la fracción de datos que han sido clasificados correctamente, por lo tanto su valor varía entre 0 y 1, y si la red funciona correctamente se espera que vaya aumentando, ya que a mayor precisión, significa que la red predice mejor. En cuanto a la función de pérdida, como se ha explicado en el capítulo 4.3.3, indica cómo de bien la red clasifica correctamente la entrada y por lo tanto se quiere minimizar. Entonces se espera que este valor vaya descendiendo conforme avanzan las épocas de entrenamiento, con el objetivo de aproximarse a un valor lo más cercano a 0. A la hora de entrenar la red, se ha indicado que se reserve un 10 % de los datos para test. Por lo tanto, en las gráficas hay información sobre el entrenamiento (en color rojo) y sobre el test (en color verde).

A continuación se va a mostrar estas gráficas para los tres modelos de redes neuronales. En las gráficas situadas a la izquierda se representa la precisión de entrenamiento y de test, y en las de la derecha el valor de la función de pérdida de entrenamiento y de test.

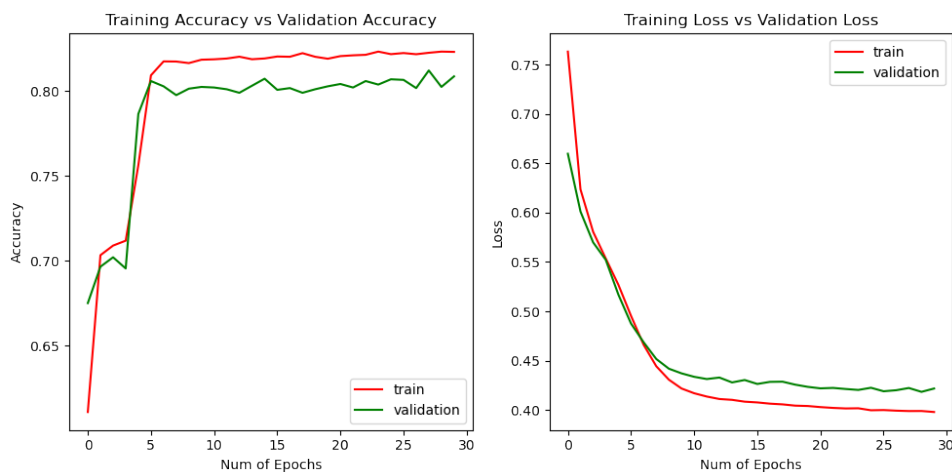


Figura 4.5: Gráficas de precisión y función de pérdida para el favour en modo max_cratio.

Como podemos observar en la figura 4.5, tanto la precisión de entrenamiento como la de test aumenta hasta estabilizarse en un 80 %. Estos valores indican una gran precisión de la red. En cuanto a la función de pérdida, tanto la pérdida de entrenamiento como la de test disminuye hasta estabilizarse en 0.4. Estos valores de pérdida son muy buenos e indican que la red predice muy bien.

Como podemos observar en la figura 4.6, tanto la precisión de entrenamiento como la de test aumenta hasta estabilizarse en un 70 %. Estos valores indican una gran precisión de la red.

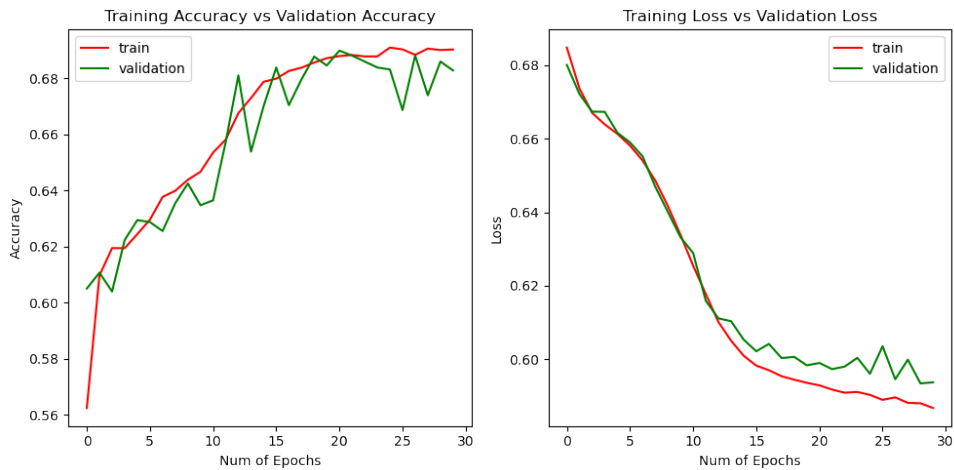


Figura 4.6: Gráficas de precisión y función de pérdida para el favour en modo max_speed.

En cuanto a la función de pérdida, tanto la pérdida de entrenamiento como la de test disminuye hasta estabilizarse en 0.6. Estos valores de pérdida son muy buenos e indican que la red predice muy bien.

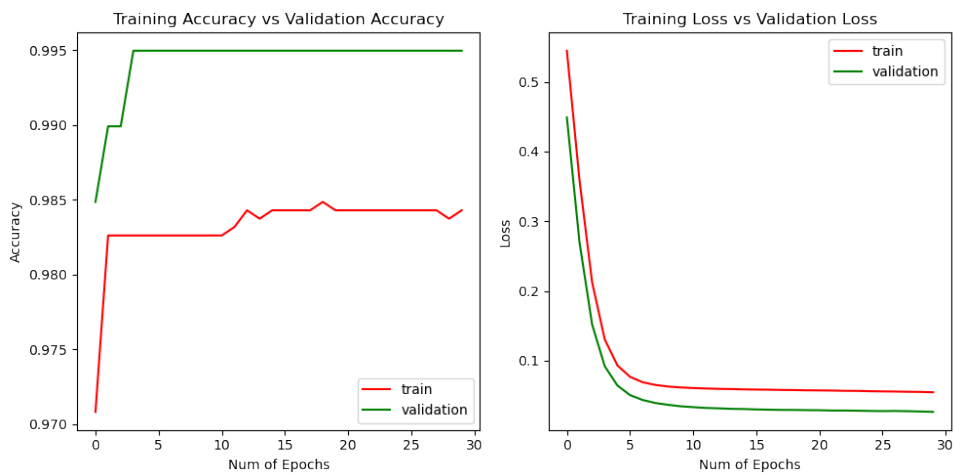


Figura 4.7: Gráficas de precisión y función de pérdida para el favour en modo balance.

Como podemos observar en la figura 4.7, tanto la precisión de entrenamiento como la de test aumenta hasta estabilizarse en prácticamente un 100%. Estos valores indican una gran precisión de la red. En cuanto a la función de pérdida, tanto la pérdida de entrenamiento como la de test disminuye hasta estabilizarse en prácticamente 0. Estos valores de pérdida son muy

buenos e indican que la red predice muy bien.

4.3.5. Estudio de las observaciones mal clasificadas

Una vez mostrados los resultados predichos por la red, podemos comparar estas predicciones con sus categorías correctas. En las siguientes imágenes podemos observar, para cada uno de los tres modelos de red neuronal, en primer lugar las predicciones de la red para todas las categorías (imágenes 4.8, 4.12 y 4.15) y en segundo lugar una comparación de cada categoría predicha con los valores predichos erróneamente de esa categoría, incluyendo un razonamiento de a qué se deben esos errores (imágenes 4.9, 4.10 y 4.11 para el modelo de red con el favour max_cratio, imágenes 4.13 y 4.13 para el modelo de red con el favour max_speed y imágenes 4.16 y 4.16 para el modelo de red con el favour balance).

En las tablas 4.2, 4.3 y 4.4 podemos observar las equivalencias de los números de las gráficas con sus categorías correspondientes. Además, las categorías que aparecen en las gráficas con un número negativo, representan la categoría con número positivo pero predichas erróneamente por la red, es decir, los errores de predicción de la red para esa categoría específica.

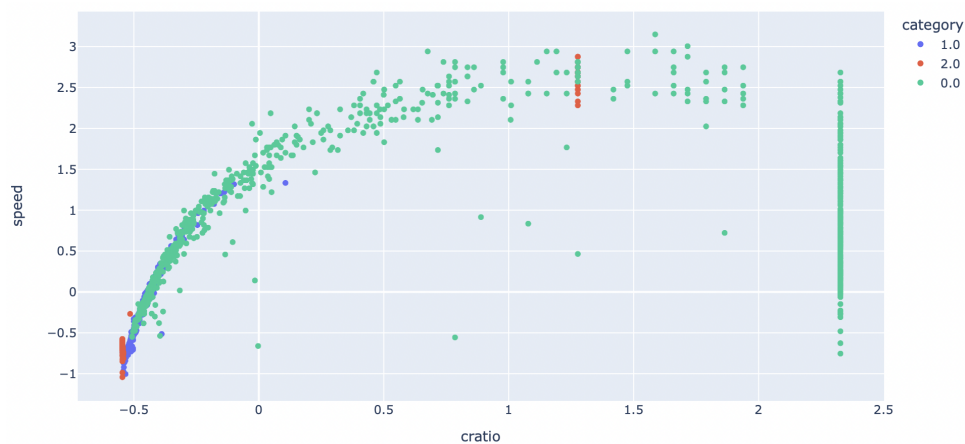


Figura 4.8: Gráfica de predicciones de la red para el favour en modo max_cratio.

En las figuras 4.18, 4.19 y 4.20 podemos observar una tabla cruzada, una técnica para el estudio de la relación existente entre dos variables categóricas, en nuestro caso entre las categorías predichas y las reales (etiquetas). En el eje horizontal se encuentran las categorías predichas por la red para cada observación del dataset, y en el eje vertical las etiquetas asociadas a esas mismas observaciones. Esta tabla cruzada se puede expresar con una matriz P de J filas y J columnas, $P = [[p_{1,1}, p_{1,2}, \dots, p_{1,J}], [p_{2,1}, p_{2,2}, \dots, p_{2,J}], \dots, [p_{J,1}, p_{J,2}, \dots, p_{J,J}]]$. Cada elemento $p_{i,j}$ representa el número de veces que la red ha predicho la categoría j cuando realmente su

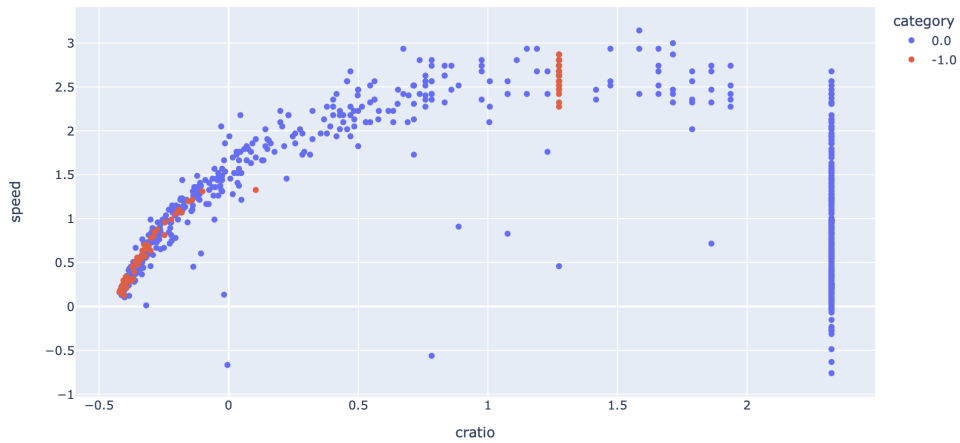


Figura 4.9: Gráfica de predicciones de la red para el favour en modo max_cratio y la categoría 0.

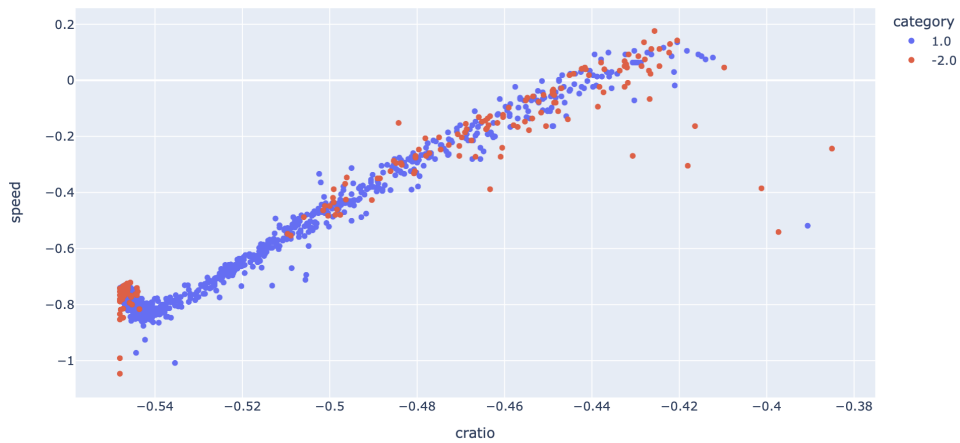


Figura 4.10: Gráfica de predicciones de la red para el favour en modo max_cratio y la categoría 1.

categoría real o etiqueta es la clase i , donde $i, j \in 0, 1, \dots, J - 1$, ya que J es el número de clases. En cada figura se representa la tabla cruzada de un modelo de red neuronal donde se optimiza cada uno de los tipos de favour.

Una vez observados los errores de las predicciones de la red para cada categoría en cada favour, el siguiente paso es analizar estos errores y estudiar por qué se han producido. Este trabajo queda planteado como línea futura de investigación.

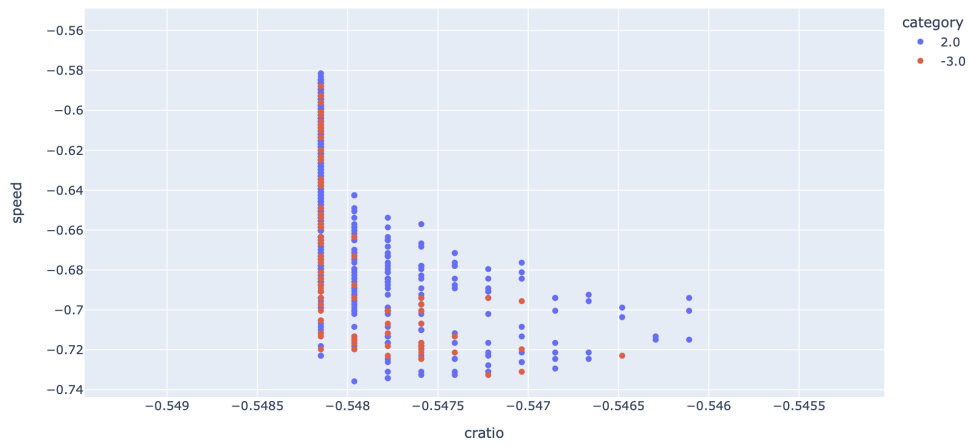


Figura 4.11: Gráfica de predicciones de la red para el favour en modo max_cratio y la categoría 2.

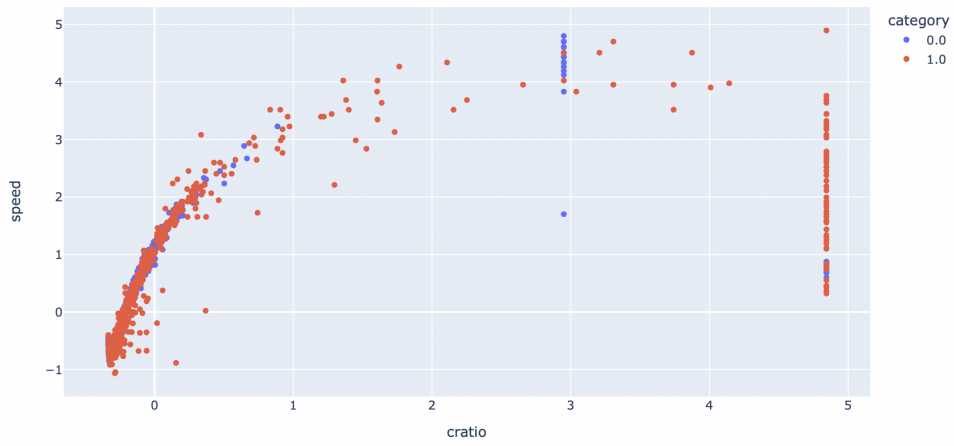


Figura 4.12: Gráfica de predicciones de la red para el favour en modo max_speed.

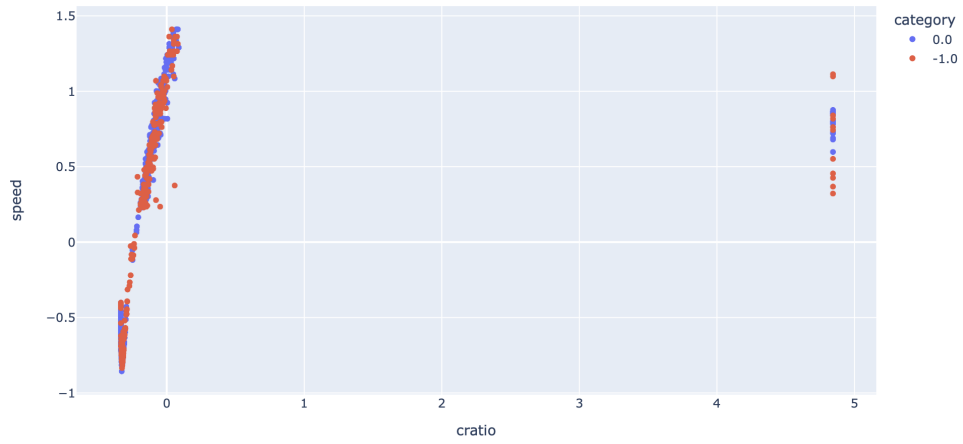


Figura 4.13: Gráfica de predicciones de la red para el favour en modo max_speed y la categoría 0.

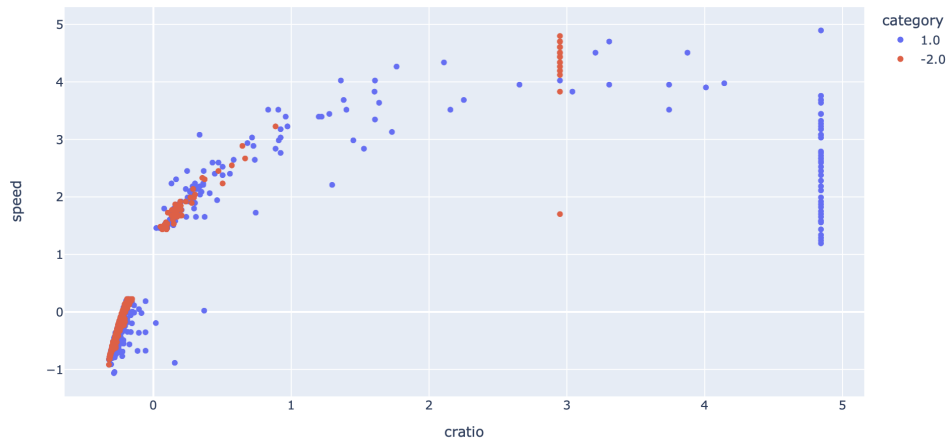


Figura 4.14: Gráfica de predicciones de la red para el favour en modo max_speed y la categoría 1.

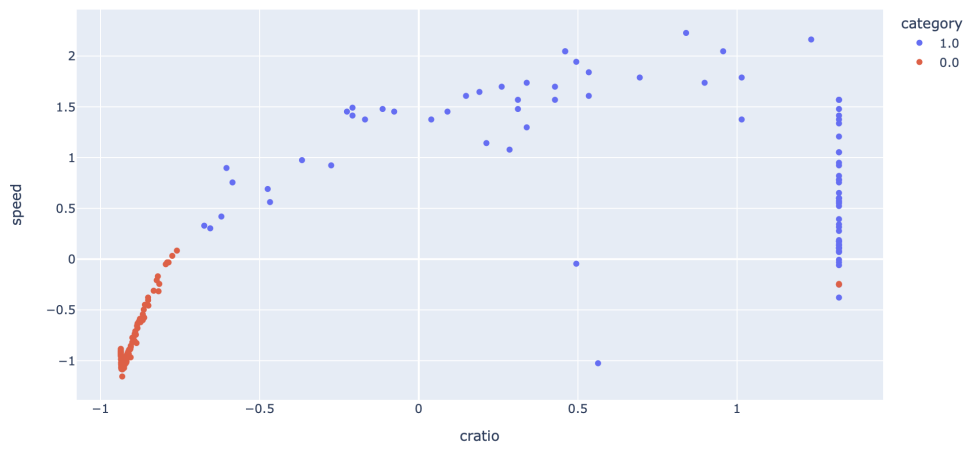


Figura 4.15: Gráfica de predicciones de la red para el favour en modo balance.

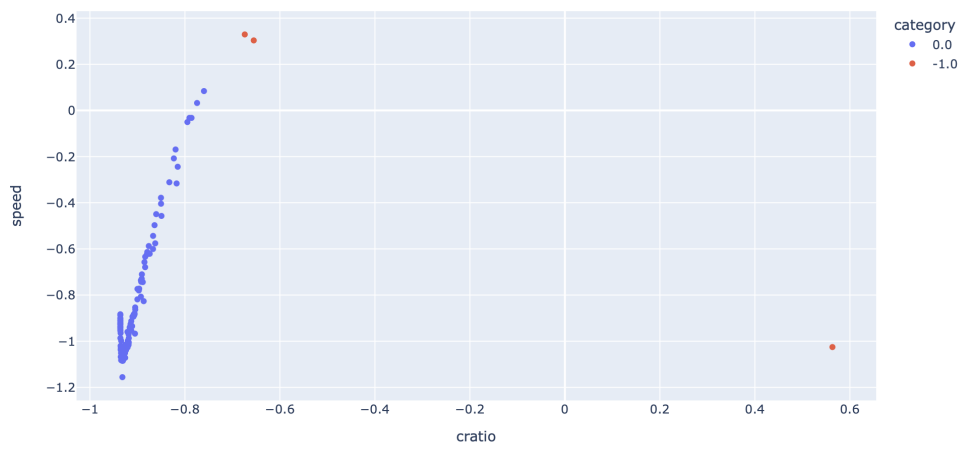


Figura 4.16: Gráfica de predicciones de la red para el favour en modo balance y la categoría 0.

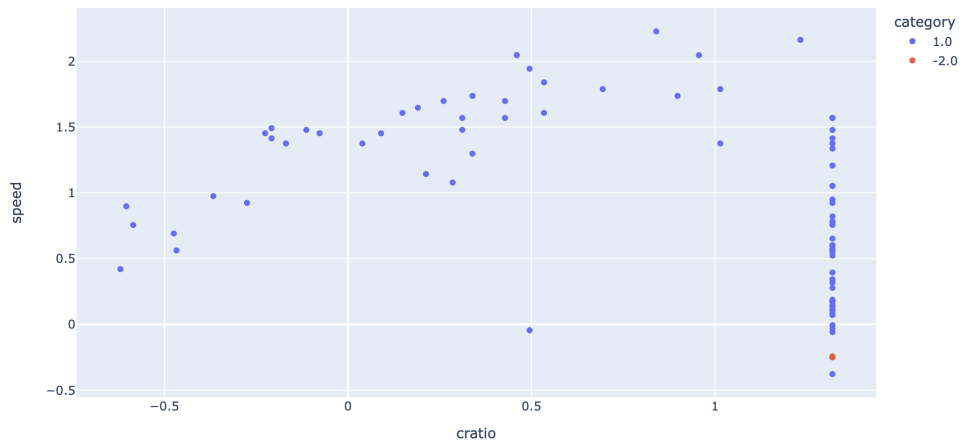


Figura 4.17: Gráfica de predicciones de la red para el favour en modo balance y la categoría 1.

col_0	0	1	2
row_0			
0	782	57	53
1	123	834	52
2	0	84	901

Figura 4.18: Tabla cruzada para el favour en modo max_cratio.

col_0	0	1
row_0		
0	929	379
1	511	1012

Figura 4.19: Tabla cruzada para el favour en modo max_speed.

col_0	0	1
row_0		
0	107	3
1	2	86

Figura 4.20: Tabla cruzada para el favour en modo balance.

Capítulo 5

Conclusiones

El objetivo principal del proyecto era desarrollar una red neuronal artificial que fuese capaz de determinar la combinación de parámetros de compresión óptima dado un dataset determinado.

Para ello se ha desarrollado un proyecto de investigación que ha implicado desde la generación de los datos, pasando por la manipulación de estos, hasta la construcción de los modelos de redes neuronales artificiales y la evaluación de su rendimiento. Para poder llevar a cabo estas tareas ha sido necesario un estudio teórico previo de las redes neuronales de clasificación, así como su aplicación práctica.

Bibliografía

- [1] Augustin Cauchy. Méthode générale pour la résolution des systemes d'équations simultanées. *Comp. Rend. Sci. Paris*, :, 25:536–538, 1847. <https://cs.uwaterloo.ca/~y328yu/classics/cauchy-en.pdf>.
- [2] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. Sgd: General analysis and improved rates. In *International Conference on Machine Learning*, pages 5200–5209. PMLR, 2019.
- [3] Geoffrey Hinton. Lecture 6e rmsprop: Divide the gradient by a running average of its recent magnitude. http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf.
- [4] B Lei P. Diederik, J. Kingma. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings. San Diego, USA, 7-9 Mayo*. <https://arxiv.org/pdf/1412.6980.pdf>.
- [5] Marta Aicart Pauner. DL-compression. <https://github.com/MartaAicart/DL-Compression.git>.
- [6] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- [7] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [8] The Blosc Development Team. C-blosc2: A fast, compressed and persistent data store library for c. <https://c-blosc2.readthedocs.io/en/latest/>.

Anexo A

Generador de datos

A.1. Código `entropy_probe.c`

En este anexo se presenta el código realizado en C que genera un documento csv donde para cada porción del dataset, se obtiene información sobre su ratio y velocidad de compresión, entre otros. El fichero del código se llama `entropy_probe.c` que incluye el programa principal y las funciones auxiliares que utiliza.

```
1  /*
2  Copyright (C) 2021 The Blosc Developers <blosc@blosc.org>
3  https://blosc.org
4  License: BSD 3-Clause (see LICENSE.txt)
5
6  Example program demonstrating use of the Blosc entropy probe codec from C
   code.
7
8  To compile this program:
9
10 $ gcc entropy_probe.c -o entropy_probe -lblosc2
11
12 To run:
13
14 $ ./entropy_probe
15
16 To run in entropy mode:
17
18 $ ./entropy_probe -e output.csv
19 */
20
21 #include <stdio.h>
```

```

22 #include "blosc2.h"
23
24 #define ENTROPY_PROBE_ID 244
25
26 #define KB 1024.
27 #define MB (1024*KB)
28 #define GB (1024*MB)
29
30 #define MAX_COPY 32U
31 #define MAX_DISTANCE 8191
32 #define MAX_FARDISTANCE (65535 + MAX_DISTANCE - 1)
33
34 // The hash length (1 << HASH_LOG2) can be tuned for performance (12 -> 15)
35 #define HASH_LOG2 (12U)
36 // #define HASH_LOG2 (13U)
37 // #define HASH_LOG2 (14U)
38
39 #define HASH_FUNCTION(v, s, h) { \
40     v = (s * 2654435761U) >> (32U - h); \
41 }
42
43 #define BLOSC_LZ_READU16(p) *((const uint16_t*)(p))
44 #define BLOSC_LZ_READU32(p) *((const uint32_t*)(p))
45
46 #define LITERAL2(ip, anchor, copy) { \
47     oc++; anchor++; \
48     ip = anchor; \
49     copy++; \
50     if (copy == MAX_COPY) { \
51         copy = 0; \
52         oc++; \
53     } \
54 }
55
56
57 static uint8_t *get_run(uint8_t *ip, const uint8_t *ip_bound, const uint8_t *
    ref) {
58     uint8_t x = ip[-1];
59     int64_t value, value2;
60     /* Broadcast the value for every byte in a 64-bit register */
61     memset(&value, x, 8);
62     /* safe because the outer check against ip limit */
63     while (ip < (ip_bound - sizeof(int64_t))) {
64         value2 = ((int64_t *) ref)[0];
65         if (value != value2) {
66             /* Return the byte that starts to differ */
67             while (*ref++ == x) ip++;
68             return ip;
69         } else {
70             ip += 8;
71             ref += 8;
72         }

```

```

73     }
74     /* Look into the remainder */
75     while ((ip < ip_bound) && (*ref++ == x)) ip++;
76     return ip;
77 }
78
79
80 static uint8_t *get_match(uint8_t *ip, const uint8_t *ip_bound, const uint8_t *
      ref) {
81     while (ip < (ip_bound - sizeof(int64_t))) {
82         if (*(int64_t *) ref != *(int64_t *) ip) {
83             /* Return the byte that starts to differ */
84             while (*ref++ == *ip++) {}
85             return ip;
86         } else {
87             ip += sizeof(int64_t);
88             ref += sizeof(int64_t);
89         }
90     }
91     /* Look into the remainder */
92     while ((ip < ip_bound) && (*ref++ == *ip++)) {}
93     return ip;
94 }
95
96
97 static uint8_t *get_run_or_match(uint8_t *ip, uint8_t *ip_bound, const uint8_t
      *ref, bool run) {
98     if (run) {
99         ip = get_run(ip, ip_bound, ref);
100    } else {
101        ip = get_match(ip, ip_bound, ref);
102    }
103
104    return ip;
105 }
106
107
108 // Get a guess for the compressed size of a buffer
109 static float get_cratio(const uint8_t *ibase, int maxlen, int minlen, int
      ipshift) {
110     const uint8_t *ip = ibase;
111     int32_t oc = 0;
112     const uint16_t hashlen = (1U << (uint8_t) HASH_LOG2);
113     uint16_t htab[1U << (uint8_t) HASH_LOG2];
114     uint32_t hval;
115     uint32_t seq;
116     uint8_t copy;
117     // Make a tradeoff between testing too much and too little
118     uint16_t limit = (maxlen > hashlen) ? hashlen : maxlen;
119     const uint8_t *ip_bound = ibase + limit - 1;
120     const uint8_t *ip_limit = ibase + limit - 12;
121

```

```

122 // Initialize the hash table to distances of 0
123 memset(htab, 0, hashlen * sizeof(uint16_t));
124
125 /* we start with literal copy */
126 copy = 4;
127 oc += 5;
128
129 /* main loop */
130 while (ip < ip_limit) {
131     const uint8_t *ref;
132     unsigned distance;
133     const uint8_t *anchor = ip;    /* comparison starting-point */
134
135     /* find potential match */
136     seq = BLOSCLZ_READU32(ip);
137     HASH_FUNCTION(hval, seq, HASH_LOG2)
138     ref = ibase + htab[hval];
139
140     /* calculate distance to the match */
141     distance = (unsigned int) (anchor - ref);
142
143     /* update hash table */
144     htab[hval] = (uint16_t) (anchor - ibase);
145
146     if (distance == 0 || (distance >= MAX_FARDISTANCE)) {
147         LITERAL2(ip, anchor, copy)
148         continue;
149     }
150
151     /* is this a match? check the first 4 bytes */
152     if (BLOSCLZ_READU32(ref) == BLOSCLZ_READU32(ip)) {
153         ref += 4;
154     } else {
155         /* no luck, copy as a literal */
156         LITERAL2(ip, anchor, copy)
157         continue;
158     }
159
160     /* last matched byte */
161     ip = anchor + 4;
162
163     /* distance is biased */
164     distance--;
165
166     /* get runs or matches; zero distance means a run */
167     ip = get_run_or_match((uint8_t*)ip, (uint8_t*)ip_bound, ref, !distance);
168
169     ip -= ipshift;
170     int32_t len = (int32_t) (ip - anchor);
171     if (len < minlen) {
172         LITERAL2(ip, anchor, copy)
173         continue;

```

```

174     }
175
176     /* if we haven't copied anything, adjust the output counter */
177     if (!copy)
178         oc--;
179     /* reset literal counter */
180     copy = 0;
181
182     /* encode the match */
183     if (distance < MAX_DISTANCE) {
184         if (len >= 7) {
185             oc += ((len - 7) / 255) + 1;
186         }
187         oc += 2;
188     } else {
189         /* far away, but not yet in the another galaxy... */
190         if (len >= 7) {
191             oc += ((len - 7) / 255) + 1;
192         }
193         oc += 4;
194     }
195
196     /* update the hash at match boundary */
197     seq = BLOSCCLZ_READU32(ip);
198     HASH_FUNCTION(hval, seq, HASH_LOG2)
199     htab[hval] = (uint16_t) (ip++ - ibase);
200     ip++;
201     /* assuming literal copy */
202     oc++;
203 }
204
205 float ic = (float) (ip - ibase);
206 return ic / (float) oc;
207 }
208
209
210 int entropy_probe(const uint8_t *input, int32_t input_len,
211                  uint8_t *output, int32_t output_len,
212                  uint8_t meta,
213                  blosc2_cparams *cparams, const void *chunk) {
214     if (cparams->splitmode != BLOSC_ALWAYS_SPLIT) {
215         BLOSC_TRACE_ERROR("Entropy probe can only be used in SPLIT mode. Aborting
216         !");
217         return BLOSC2_ERROR_CODE_PARAM;
218     }
219     // Get the cratio. minlen and ipshift are decent defaults, but one can try
220     // with (4, 4) or (3, 4) or (4, 3)
221     float cratio = get_cratio(input, input_len, 3, 3);
222     int cbytes = (int) ((float)input_len / cratio);
223     if (cbytes > input_len) {
224         cbytes = input_len;
225     }

```

```

224     return cbytes;
225 }
226
227
228 // This function receives an instrumented chunk having nstreams
229 int extr_data(FILE *csv_file, int nchunk, uint8_t *chunk, int nstreams, int
    category) {
230     blosc2_instr *instr_data = (blosc2_instr *) chunk;
231
232     for (int nstream = 0; nstream < nstreams; nstream++) {
233         float cratio = instr_data->cratio;
234         float cspeed = instr_data->cspeed;
235         bool special_val = instr_data->flags[0];
236         if (!special_val) {
237             // Fill csv file
238             int special_vals = 0;
239             fprintf(csv_file, "%.3g, %.3g, %d, %d, %d\n", cratio, cspeed,
                special_vals, nchunk, category);
240             printf("Chunk %d, block %d: cratio %.3g, speed %.3g\n", nchunk, nstream,
                cratio, cspeed);
241         } else {
242             // Fill csv file
243             int special_vals = 1;
244             fprintf(csv_file, "%.3g, %.3g, %d, %d, %d\n", cratio, cspeed,
                special_vals, nchunk, category);
245             printf("Chunk %d, block %d: cratio %.3g, speed %.3g\n", nchunk, nstream
                , cratio, cspeed);
246         }
247         instr_data++;
248     }
249
250     return 0;
251 }
252
253 //Categorize the data according to the combination of filter-codec-splitmode
254 int categorize_data(char codec, char filter, char splitmode) {
255     printf(" AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA el nombre del codec es: %c\n", codec
    );
256
257     int category;
258     if (codec == BLOSC_BLOSC LZ && filter == BLOSC_NOFILTER && splitmode ==
    BLOSC_ALWAYS_SPLIT) {
259         category = 0;
260     } else if (codec == BLOSC_BLOSC LZ && filter == BLOSC_SHUFFLE && splitmode
    == BLOSC_ALWAYS_SPLIT) {
261         category = 1;
262     } else if (codec == BLOSC_BLOSC LZ && filter == BLOSC_BITSHUFFLE &&
    splitmode == BLOSC_ALWAYS_SPLIT) {
263         category = 2;
264     } else if (codec == BLOSC_LZ4 && filter == BLOSC_NOFILTER && splitmode ==
    BLOSC_ALWAYS_SPLIT) {
265         category = 3;

```



```

266     } else if (codec == BLOSC_LZ4 && filter == BLOSC_SHUFFLE && splitmode ==
BLOSC_ALWAYS_SPLIT) {
267         category = 4;
268     } else if (codec == BLOSC_LZ4 && filter == BLOSC_BITSHUFFLE && splitmode ==
BLOSC_ALWAYS_SPLIT) {
269         category = 5;
270     } else if (codec == BLOSC_LZ4HC && filter == BLOSC_NOFILTER && splitmode ==
BLOSC_ALWAYS_SPLIT) {
271         category = 6;
272     } else if (codec == BLOSC_LZ4HC && filter == BLOSC_SHUFFLE && splitmode ==
BLOSC_ALWAYS_SPLIT) {
273         category = 7;
274     } else if (codec == BLOSC_LZ4HC && filter == BLOSC_BITSHUFFLE && splitmode
== BLOSC_ALWAYS_SPLIT) {
275         category = 8;
276     } else if (codec == BLOSC_ZLIB && filter == BLOSC_NOFILTER && splitmode ==
BLOSC_ALWAYS_SPLIT) {
277         category = 9;
278     } else if (codec == BLOSC_ZLIB && filter == BLOSC_SHUFFLE && splitmode ==
BLOSC_ALWAYS_SPLIT) {
279         category = 10;
280     } else if (codec == BLOSC_ZLIB && filter == BLOSC_BITSHUFFLE && splitmode
== BLOSC_ALWAYS_SPLIT) {
281         category = 11;
282     } else if (codec == BLOSC_ZSTD && filter == BLOSC_NOFILTER && splitmode ==
BLOSC_ALWAYS_SPLIT) {
283         category = 12;
284     } else if (codec == BLOSC_ZSTD && filter == BLOSC_SHUFFLE && splitmode ==
BLOSC_ALWAYS_SPLIT) {
285         category = 13;
286     } else if (codec == BLOSC_ZSTD && filter == BLOSC_BITSHUFFLE && splitmode
== BLOSC_ALWAYS_SPLIT) {
287         category = 14;
288
289     } else if (codec == BLOSC_BLOSC LZ && filter == BLOSC_NOFILTER && splitmode
== BLOSC_NEVER_SPLIT) {
290         category = 15;
291     } else if (codec == BLOSC_BLOSC LZ && filter == BLOSC_SHUFFLE && splitmode
== BLOSC_NEVER_SPLIT) {
292         category = 16;
293     } else if (codec == BLOSC_BLOSC LZ && filter == BLOSC_BITSHUFFLE &&
splitmode == BLOSC_NEVER_SPLIT) {
294         category = 17;
295     } else if (codec == BLOSC_LZ4 && filter == BLOSC_NOFILTER && splitmode ==
BLOSC_NEVER_SPLIT) {
296         category = 18;
297     } else if (codec == BLOSC_LZ4 && filter == BLOSC_SHUFFLE && splitmode ==
BLOSC_NEVER_SPLIT) {
298         category = 19;
299     } else if (codec == BLOSC_LZ4 && filter == BLOSC_BITSHUFFLE && splitmode ==
BLOSC_NEVER_SPLIT) {
300         category = 20;

```

```

301     } else if (codec == BLOSC_LZ4HC && filter == BLOSC_NOFILTER && splitmode ==
302         BLOSC_NEVER_SPLIT) {
303         category = 21;
304     } else if (codec == BLOSC_LZ4HC && filter == BLOSC_SHUFFLE && splitmode ==
305         BLOSC_NEVER_SPLIT) {
306         category = 22;
307     } else if (codec == BLOSC_LZ4HC && filter == BLOSC_BITSHUFFLE && splitmode
308         == BLOSC_NEVER_SPLIT) {
309         category = 23;
310     } else if (codec == BLOSC_ZLIB && filter == BLOSC_NOFILTER && splitmode ==
311         BLOSC_NEVER_SPLIT) {
312         category = 24;
313     } else if (codec == BLOSC_ZLIB && filter == BLOSC_SHUFFLE && splitmode ==
314         BLOSC_NEVER_SPLIT) {
315         category = 25;
316     } else if (codec == BLOSC_ZLIB && filter == BLOSC_BITSHUFFLE && splitmode
317         == BLOSC_NEVER_SPLIT) {
318         category = 26;
319     } else if (codec == BLOSC_ZSTD && filter == BLOSC_NOFILTER && splitmode ==
320         BLOSC_NEVER_SPLIT) {
321         category = 27;
322     } else if (codec == BLOSC_ZSTD && filter == BLOSC_SHUFFLE && splitmode ==
323         BLOSC_NEVER_SPLIT) {
324         category = 28;
325     } else if (codec == BLOSC_ZSTD && filter == BLOSC_BITSHUFFLE && splitmode
326         == BLOSC_NEVER_SPLIT) {
327         category = 29;
328     }
329
330     return category;
331 }
332
333 void print_compress_info(void) {
334     char *name = NULL, *version = NULL;
335     int ret;
336
337     printf("Blosc version: %s (%s)\n", BLOSC_VERSION_STRING, BLOSC_VERSION_DATE);
338
339     printf("List of supported compressors in this build: %s\n",
340         blosc_list_compressors());
341
342     printf("Supported compression libraries:\n");
343     ret = blosc_get_complib_info("blosclz", &name, &version);
344     if (ret >= 0) printf("  %s: %s\n", name, version);
345     free(name);
346     free(version);
347     ret = blosc_get_complib_info("lz4", &name, &version);
348     if (ret >= 0) printf("  %s: %s\n", name, version);
349     free(name);
350     free(version);
351     ret = blosc_get_complib_info("zlib", &name, &version);
352     if (ret >= 0) printf("  %s: %s\n", name, version);

```

```

344     free(name);
345     free(version);
346     ret = blosc_get_complib_info("zstd", &name, &version);
347     if (ret >= 0) printf("  %s: %s\n", name, version);
348     free(name);
349     free(version);
350 }
351
352
353 int main(int argc, char *argv[]) {
354     char usage[256];
355     char csv_filename[256];
356     char data_filename[256];
357     bool entropy_probe_mode = false;
358
359     print_compress_info();
360
361     strcpy(usage, "Usage: entropy_probe [-e] data_filename");
362
363     if (argc < 2) {
364         printf("%s\n", usage);
365         exit(1);
366     }
367
368     if (argc >= 3) {
369         if (strcmp("-e", argv[1]) != 0) {
370             printf("%s\n", usage);
371             exit(1);
372         }
373         strcpy(data_filename, argv[2]);
374         entropy_probe_mode = true;
375     }
376     else {
377         strcpy(data_filename, argv[1]);
378     }
379     printf("fitxer %s\n", data_filename);
380     blosc2_cparams cparams = BLOSC2_CPARAMS_DEFAULTS;
381     cparams.instr_codec = true;
382     cparams.blocksize = 256 * (int)KB;
383     cparams.splitmode = BLOSC_ALWAYS_SPLIT;
384
385     if (entropy_probe_mode) {
386         // The entropy probe detector is meant to always be used in SPLIT mode
387         cparams.splitmode = BLOSC_ALWAYS_SPLIT;
388     }
389     cparams.typesize = sizeof(float);
390
391     if (entropy_probe_mode) {
392         blosc2_codec udcodec;
393         udcodec.compcode = ENTROPY_PROBE_ID;
394         udcodec.compver = 1;
395         udcodec.complib = 1;

```

```

396     udcodec.compname = "entropy_probe";
397     udcodec.encoder = entropy_probe;
398     udcodec.decoder = NULL;
399     blosc2_register_codec(&udcodec);
400     cparams.compcode = ENTROPY_PROBE_ID;
401 } //else {
402 //     cparams.compcode = BLOSC_BLOSC LZ;
403 //}
404
405 blosc2_dparams dparams = BLOSC2_DPARAMS_DEFAULTS;
406
407 blosc2_schunk *schunk = blosc2_schunk_open(data_filename);
408 if (schunk == NULL) {
409     printf("Cannot open the data file\n");
410     exit(1);
411 }
412 uint8_t *chunk = malloc(schunk->chunksz);
413 uint8_t *chunk2 = malloc(schunk->chunksz);
414 printf("nchunks in dataset: %d\n", schunk->nchunks);
415
416 if (!entropy_probe_mode) {
417     // Loop over different filters and codecs
418     int v_codecs[] = {0, 1, 2, 4, 5};
419     int v_splits[] = {1, 2};
420
421     for (int ncodec = 0; ncodec < sizeof(v_codecs); ++ncodec) {
422         for (int nfilter = 0; nfilter <= BLOSC_BITSHUFFLE; ++nfilter) {
423             for (int nsplit = 0; nsplit < BLOSC_NEVER_SPLIT; ++nsplit) {
424
425                 cparams.splitmode = v_splits[nsplit];
426                 cparams.compcode = v_codecs[ncodec];
427                 cparams.filters[BLOSC2_MAX_FILTERS - 1] = nfilter;
428                 blosc2_context *cctx = blosc2_create_cctx(cparams);
429                 blosc2_context *dctx = blosc2_create_dctx(dparams);
430
431                 const char *compname;
432                 switch (cparams.compcode) {
433                     case BLOSC_BLOSC LZ:
434                         compname = BLOSC_BLOSC LZ_COMPNAME;
435                         break;
436                     case BLOSC_LZ4:
437                         compname = BLOSC_LZ4_COMPNAME;
438                         break;
439                     case BLOSC_LZ4HC:
440                         compname = BLOSC_LZ4HC_COMPNAME;
441                         break;
442                     case BLOSC_ZLIB:
443                         compname = BLOSC_ZLIB_COMPNAME;
444                         break;
445                     case BLOSC_ZSTD:
446                         compname = BLOSC_ZSTD_COMPNAME;
447                         break;

```

```

448         case ENTROPY_PROBE_ID:
449             compname = "entropy";
450             break;
451         default:
452             printf("Unsupported codec!");
453             exit(1);
454     }
455
456     char *sfilter;
457     switch (nfilter) {
458         case BLOSC_NOFILTER:
459             sfilter = "nofilter";
460             break;
461         case BLOSC_SHUFFLE:
462             sfilter = "shuffle";
463             break;
464         case BLOSC_BITSHUFFLE:
465             sfilter = "bitshuffle";
466             break;
467         default:
468             printf("Unsupported filter!");
469             exit(1);
470     }
471
472     char *ssplit;
473     switch (cparams.splitmode) {
474         case BLOSC_ALWAYS_SPLIT:
475             ssplit = "split";
476             break;
477         case BLOSC_NEVER_SPLIT:
478             ssplit = "nosplit";
479             break;
480
481         default:
482             printf("Unsupported splitmode!");
483             exit(1);
484     }
485
486     // Create csv file
487     sprintf(csv_filename, "%s-%s-%s.csv", compname, sfilter,
ssplit);
488
489     printf("Codec : %s\n", compname);
490     printf("Filter : %s\n", sfilter);
491     printf("Splitmode: %s\n", ssplit);
492
493     FILE *csv_file = fopen(csv_filename, "w");
494     if (csv_file == NULL) {
495         printf("Error creating the file\n");
496         return -1;
497     }
498     int category;

```

```

499         if (strcmp(compname, BLOSC_BLOSC_LZ_COMPNAME) == 0 && strcmp(
sfilter, "nofilter") == 0 && strcmp(ssplit, "split") == 0) {
500             category = 0;
501         } else if (strcmp(compname, BLOSC_BLOSC_LZ_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "split") == 0) {
502             category = 1;
503         } else if (strcmp(compname, BLOSC_BLOSC_LZ_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "split") == 0) {
504             category = 2;
505         } else if (strcmp(compname, BLOSC_LZ4_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "split") == 0) {
506             category = 3;
507         } else if (strcmp(compname, BLOSC_LZ4_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "split") == 0) {
508             category = 4;
509         } else if (strcmp(compname, BLOSC_LZ4_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "split") == 0) {
510             category = 5;
511         } else if (strcmp(compname, BLOSC_LZ4HC_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "split") == 0) {
512             category = 6;
513         } else if (strcmp(compname, BLOSC_LZ4HC_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "split") == 0) {
514             category = 7;
515         } else if (strcmp(compname, BLOSC_LZ4HC_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "split") == 0) {
516             category = 8;
517         } else if (strcmp(compname, BLOSC_ZLIB_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "split") == 0) {
518             category = 9;
519         } else if (strcmp(compname, BLOSC_ZLIB_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "split") == 0) {
520             category = 10;
521         } else if (strcmp(compname, BLOSC_ZLIB_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "split") == 0) {
522             category = 11;
523         } else if (strcmp(compname, BLOSC_ZSTD_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "split") == 0) {
524             category = 12;
525         } else if (strcmp(compname, BLOSC_ZSTD_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "split") == 0) {
526             category = 13;
527         } else if (strcmp(compname, BLOSC_ZSTD_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "split") == 0) {
528             category = 14;
529
530         } else if (strcmp(compname, BLOSC_BLOSC_LZ_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "nosplit") == 0) {
531             category = 15;
532         } else if (strcmp(compname, BLOSC_BLOSC_LZ_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
533             category = 16;

```

```

534         } else if (strcmp(compname, BLOSC_BLOSC_LZ_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
535             category = 17;
536         } else if (strcmp(compname, BLOSC_LZ4_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "nosplit") == 0) {
537             category = 18;
538         } else if (strcmp(compname, BLOSC_LZ4_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
539             category = 19;
540         } else if (strcmp(compname, BLOSC_LZ4_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
541             category = 20;
542         } else if (strcmp(compname, BLOSC_LZ4HC_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "nosplit") == 0) {
543             category = 21;
544         } else if (strcmp(compname, BLOSC_LZ4HC_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
545             category = 22;
546         } else if (strcmp(compname, BLOSC_LZ4HC_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
547             category = 23;
548         } else if (strcmp(compname, BLOSC_ZLIB_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "nosplit") == 0) {
549             category = 24;
550         } else if (strcmp(compname, BLOSC_ZLIB_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
551             category = 25;
552         } else if (strcmp(compname, BLOSC_ZLIB_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
553             category = 26;
554         } else if (strcmp(compname, BLOSC_ZSTD_COMPNAME) == 0 &&
strcmp(sfilter, "nofilter") == 0 && strcmp(ssplit, "nosplit") == 0) {
555             category = 27;
556         } else if (strcmp(compname, BLOSC_ZSTD_COMPNAME) == 0 &&
strcmp(sfilter, "shuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
557             category = 28;
558         } else if (strcmp(compname, BLOSC_ZSTD_COMPNAME) == 0 &&
strcmp(sfilter, "bitshuffle") == 0 && strcmp(ssplit, "nosplit") == 0) {
559             category = 29;
560         }
561
562         fprintf(csv_file, "cratio, speed, special_vals, nchunk,
category\n");
563         for (int nchunk = 0; nchunk < schunk->nchunks; nchunk++) {
564             printf("decompressing chunk # %d (out of %d)\n", nchunk,
schunk->nchunks);
565             int dsize = blosc2_schunk_decompress_chunk(schunk, nchunk
, chunk, schunk->chunksize);
566             if (dsize < 0) {
567                 printf("Error decompressing chunk in schunk. Error
code: %d\n", dsize);
568                 return dsize;

```

```

569         }
570         int csize = blosc2_compress_ctx(cctx, chunk, dsize,
chunk2, schunk->chunksize);
571         if (csize < 0) {
572             printf("Error compressing chunk. Error code: %d\n",
csize);
573             return csize;
574         }
575         int dsize2 = blosc2_decompress_ctx(dctx, chunk2, csize,
chunk, dsize);
576         if (dsize2 < 0) {
577             printf("Error decompressing chunk. Error code: %d\n
", dsize2);
578             return dsize2;
579         }
580
581         int nstreams = dsize2 / (int) sizeof(blosc2_instr);
582         printf("Chunk %d data with %d streams:\n", nchunk,
nstreams);
583
584         int categoria = categorize_data(compname, *sfilter, *
ssplit);
585         extr_data(csv_file, nchunk, chunk, nstreams, category);
586     }
587     fclose(csv_file);
588 }
589 }
590 }
591 }
592
593 } else {
594     // Loop over different filters
595     for (int nfilter = 0; nfilter <= BLOSC_BITSHUFFLE; ++nfilter) {
596         cparams.filters[BLOSC2_MAX_FILTERS - 1] = nfilter;
597         blosc2_context *cctx = blosc2_create_ctx(cparams);
598         blosc2_context *dctx = blosc2_create_dctx(dparams);
599
600         const char *compname;
601         switch (cparams.compcode) {
602             case BLOSC_BLOSC LZ:
603                 compname = BLOSC_BLOSC LZ_COMPNAME;
604                 break;
605             case BLOSC_LZ4:
606                 compname = BLOSC_LZ4_COMPNAME;
607                 break;
608             case BLOSC_LZ4HC:
609                 compname = BLOSC_LZ4HC_COMPNAME;
610                 break;
611             case BLOSC_ZLIB:
612                 compname = BLOSC_ZLIB_COMPNAME;
613                 break;
614             case BLOSC_ZSTD:

```



```

615         compname = BLOSC_ZSTD_COMPNAME;
616         break;
617     case ENTROPY_PROBE_ID:
618         compname = "entropy";
619         break;
620     default:
621         printf("Unsupported codec!");
622         exit(1);
623 }
624
625 char *sfilter;
626 switch (nfilter) {
627     case BLOSC_NOFILTER:
628         sfilter = "nofilter";
629         break;
630     case BLOSC_SHUFFLE:
631         sfilter = "shuffle";
632         break;
633     case BLOSC_BITSHUFFLE:
634         sfilter = "bitshuffle";
635         break;
636     default:
637         printf("Unsupported filter!");
638         exit(1);
639 }
640
641 // Create csv file
642 sprintf(csv_filename, "%s-%s.csv", compname, sfilter);
643 FILE *csv_file = fopen(csv_filename, "w");
644 if (csv_file == NULL) {
645     printf("Error creating the file\n");
646     return -1;
647 }
648
649 fprintf(csv_file, "cratio, speed, special_vals, nchunk, category\n");
650 for (int nchunk = 0; nchunk < schunk->nchunks; nchunk++) {
651     printf("decompressing chunk # %d (out of %d)\n", nchunk, schunk->
nchunks);
652     int dsize = blocs2_schunk_decompress_chunk(schunk, nchunk, chunk,
schunk->chunksize);
653     if (dsize < 0) {
654         printf("Error decompressing chunk in schunk. Error code: %d\
n", dsize);
655         return dsize;
656     }
657     int csize = blocs2_compress_ctx(cctx, chunk, dsize, chunk2,
schunk->chunksize);
658     if (csize < 0) {
659         printf("Error compressing chunk. Error code: %d\n", csize);
660         return csize;
661     }

```

```

662     int dsize2 = blosc2_decompress_ctx(dctx, chunk2, csize, chunk,
    dsize);
663     if (dsize2 < 0) {
664         printf("Error decompressing chunk. Error code: %d\n", dsize2
    );
665         return dsize2;
666     }
667
668     int nstreams = dsize2 / (int) sizeof(blosc2_instr);
669     printf("Chunk %d data with %d streams:\n", nchunk, nstreams);
670
671     int category = -1;
672     extr_data(csv_file, nchunk, chunk, nstreams, category);
673 }
674
675     fclose(csv_file);
676
677 }
678 }
679
680 /* Free resources */
681 blosc2_schunk_free(schunk);
682 free(chunk);
683 free(chunk2);
684 printf("Success!\n");
685
686 return 0;
687 }

```

Anexo B

Análisis de datos

Los siguientes códigos de las secciones B.1, B.2 y B.3 estan publicados en mi repositorio de Git Hub [5].

B.1. Código Neural-Network-CRATIO.ipynb

En este anexo se presenta el código realizado en Jupyter Notebook que prepara y manipula los datos, para finalmente construir y evaluar el modelo de red neuronal donde se optimiza el favour de max_cratio.El fichero del código se llama **Neural-Network-CRATIO.ipynb**.

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 # In[128]:
5
6
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 import plotly.express as px
10 import numpy as np
11
12
13 # In[129]:
14
15
16 from pathlib import Path
17 path = Path('../data_inputs_temp/')
18
19
```

```

20 # In[130]:
21
22
23 ficheros = ['entropy-shuffle', 'zstd-shuffle-split', 'zlib-shuffle-split', '
    lz4hc-shuffle-split']
24
25 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
26 dataframes1 = {}
27 nsplits = [55296] * 4
28
29 for i in range(0, len(ficheros)):
30     dataframes1[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
31     codec_filter = np.repeat(ficheros[i], nsplits[i])
32     dataframes1[ficheros[i]]['codec_filter'] = codec_filter
33     if i == 0:
34         categoria = np.repeat(str(30+i), nsplits[i])
35     else:
36         categoria = np.repeat(str(i-1), nsplits[i])
37     dataframes1[ficheros[i]]['categoria'] = categoria
38
39     dataframes1[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
    speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
    category'}, inplace = True)
40
41
42 # In[131]:
43
44
45 from pathlib import Path
46 path = Path('../data_inputs')
47
48
49 # In[132]:
50
51
52 ficheros = ['entropy-shuffle', 'zstd-shuffle-split', 'zlib-shuffle-split', '
    lz4hc-shuffle-split']
53
54 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
55 dataframes2 = {}
56 nsplits = [55296] * 4
57
58 for i in range(0, len(ficheros)):
59     dataframes2[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
60     codec_filter = np.repeat(ficheros[i], nsplits[i])
61     dataframes2[ficheros[i]]['codec_filter'] = codec_filter
62     if i == 0:
63         categoria = np.repeat(str(30+i), nsplits[i])
64     else:
65         categoria = np.repeat(str(i-1), nsplits[i])
66     dataframes2[ficheros[i]]['categoria'] = categoria
67

```

```

68     dataframes2[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
69
70
71 # In[133]:
72
73
74 from pathlib import Path
75 path = Path('../data_inputs_wind/')
76
77
78 # In[134]:
79
80
81 ficheros = ['entropy-shuffle', 'zstd-shuffle-split', 'zlib-shuffle-split', '
lz4hc-shuffle-split']
82
83 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
84 dataframes3 = {}
85 nsplits = [55296] * 4
86
87 for i in range(0, len(ficheros)):
88     dataframes3[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
89     codec_filter = np.repeat(ficheros[i], nsplits[i])
90     dataframes3[ficheros[i]]['codec_filter'] = codec_filter
91     if i == 0:
92         categoria = np.repeat(str(30+i), nsplits[i])
93     else:
94         categoria = np.repeat(str(i-1), nsplits[i])
95     dataframes3[ficheros[i]]['categoria'] = categoria
96
97     dataframes3[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
98
99
100 # In[135]:
101
102
103 df_fig = pd.DataFrame()
104 for i in range(0, len(ficheros)):
105     df_fig = pd.concat([df_fig, dataframes1[ficheros[i]],
106                        dataframes2[ficheros[i]], dataframes3[ficheros[i]]],
107                        axis = 0)
108
109 # In[136]:
110
111
112 df_fig = df_fig[df_fig['special_vals'] == 0]
113

```

```

114
115 # In[137]:
116
117
118 fig = px.scatter(df_fig, "cratio", "speed", color = "categoria")
119 fig.show()
120
121
122 # 0: zstd
123 #
124 # 1: zlib
125 #
126 # 2: lz4hc
127 #
128 # 30: entropy
129
130 # # Dataframe filter: SHUFFLE
131
132 # In[138]:
133
134
135 df_entropy = pd.DataFrame()
136 df_entropy = pd.concat([df_entropy, dataframes1[ficheros[0]], dataframes2[
    ficheros[0]],
137                        dataframes3[ficheros[0]]], axis = 0)
138
139 df_entropy.rename(columns = {'cratio':'cratio31', 'speed':'speed31', '
    special_vals':'special_vals31'}, inplace = True)
140 df_entropy = df_entropy.drop(['nchunk', 'category', 'codec_filter', 'categoria
    '], axis = 1)
141
142
143 # In[139]:
144
145
146 df_zstd = pd.DataFrame()
147 df_zstd = pd.concat([df_zstd, dataframes1[ficheros[1]], dataframes2[ficheros
    [1]],
148                    dataframes3[ficheros[1]]], axis = 0)
149
150 df_zstd.rename(columns = {'cratio':'cratio0', 'speed':'speed0', 'special_vals
    ':'special_vals0'}, inplace = True)
151 df_zstd = df_zstd.drop(['nchunk', 'category', 'codec_filter', 'categoria'],
    axis = 1)
152
153
154 # In[140]:
155
156
157 df_zlib = pd.DataFrame()
158 df_zlib = pd.concat([df_zlib, dataframes1[ficheros[2]], dataframes2[ficheros
    [2]],

```

```

159         dataframes3[ficheros[2]], axis = 0)
160
161 df_zlib.rename(columns = {'cratio':'cratio1', 'speed':'speed1', 'special_vals
162   ':'special_vals1'}, inplace = True)
163 df_zlib = df_zlib.drop(['nchunk', 'category', 'codec_filter', 'categoria'],
164   axis = 1)
165
166 # In[141]:
167
168 df_lz4hc = pd.DataFrame()
169 df_lz4hc = pd.concat([df_lz4hc, dataframes1[ficheros[3]], dataframes2[ficheros
170   [3]],
171   dataframes3[ficheros[3]]], axis = 0)
172 df_lz4hc.rename(columns = {'cratio':'cratio2', 'speed':'speed2', 'special_vals
173   ':'special_vals2'}, inplace = True)
174 df_lz4hc = df_lz4hc.drop(['nchunk', 'category', 'codec_filter', 'categoria'],
175   axis = 1)
176
177 # In[142]:
178
179 df_shuffle = pd.DataFrame()
180 df_shuffle = pd.concat([df_shuffle, df_zstd, df_zlib, df_lz4hc, df_entropy],
181   axis = 1)
182
183 # Eliminar special_value (que lo sea en alguna de las 4 categor as)
184
185 # In[143]:
186
187
188 def delete_final(special_vals0, special_vals1, special_vals2, special_vals31):
189     if (special_vals0 == 1 or special_vals1 == 1 or special_vals2 == 1 or
190     special_vals31 == 1):
191         return 1
192     else:
193         return 0
194
195 # In[144]:
196
197
198 df_shuffle["delete_f"] = df_shuffle.apply(lambda x: delete_final(x[2], x[5], x
199   [8], x[11]), axis=1)
200
201 # In[145]:
202

```

```

203
204 df_shuffle = df_shuffle[df_shuffle['delete_f'] == 0]
205
206
207 # In[146]:
208
209
210 df_shuffle = df_shuffle.drop(['delete_f'], axis = 1)
211
212
213 # In[147]:
214
215
216 df_input_x = df_shuffle.iloc[:, [9, 10]]
217
218
219 # C lculo del SCORE
220
221 # In[148]:
222
223
224 df_shuffle = df_shuffle.drop(['cratio31', 'speed31', 'special_vals31', '
    special_vals0', 'special_vals1', 'special_vals2'], axis = 1)
225
226
227 # In[149]:
228
229
230 dfcratios = pd.DataFrame()
231 for i in range(0, 3):
232     dfcratios = pd.concat([dfcratios, df_shuffle.iloc[:, 2*i]], axis =0)
233
234
235 # In[150]:
236
237
238 desv_cratio = dfcratios.std()[0]
239 desv_cratio
240
241
242 # In[151]:
243
244
245 mean_cratio = dfcratios.mean()[0]
246 mean_cratio
247
248
249 # In[152]:
250
251
252 dfspeeds = pd.DataFrame()
253 for i in range(0, 3):

```



```

254     dfspeeds = pd.concat([dfspeeds, df_shuffle.iloc[:, 2*i+1]], axis =0)
255
256
257 # In[153]:
258
259
260 mean_speed = dfspeeds.mean()[0]
261 mean_speed
262
263
264 # In[154]:
265
266
267 desv_speed = dfspeeds.std()[0]
268 desv_speed
269
270
271 # In[155]:
272
273
274 theta = 0.96
275
276 cratio0 = df_shuffle["cratio0"]
277 speed0 = df_shuffle["speed0"]
278 ncratio0 = (cratio0 - mean_cratio) / desv_cratio
279 nspeed0 = (speed0 - mean_speed) / desv_speed
280 df_shuffle.insert(2, 0, theta*ncratio0 + (1-theta)*nspeed0)
281
282 cratio1 = df_shuffle["cratio1"]
283 speed1 = df_shuffle["speed1"]
284 ncratio1 = (cratio1 - mean_cratio) / desv_cratio
285 nspeed1 = (speed1 - mean_speed) / desv_speed
286 df_shuffle.insert(5, 1, theta*ncratio1 + (1-theta)*nspeed1)
287
288 cratio2 = df_shuffle["cratio2"]
289 speed2 = df_shuffle["speed2"]
290 ncratio2 = (cratio2 - mean_cratio) / desv_cratio
291 nspeed2 = (speed2 - mean_speed) / desv_speed
292 df_shuffle.insert(8, 2, theta*ncratio2 + (1-theta)*nspeed2)
293
294
295 # In[156]:
296
297
298 df_pos = df_shuffle
299 df_pos = df_pos.drop(['cratio0', 'speed0', 'cratio1', 'speed1', 'cratio2', '
    speed2'], axis = 1)
300
301
302 # In[157]:
303
304

```

```

305 best_categ_shuffle = df_pos.idxmax(axis=1)
306 best_categ_shuffle = pd.DataFrame(best_categ_shuffle, columns=['best_categ'])
307
308
309 # In[158]:
310
311
312 categorias = best_categ_shuffle['best_categ'].unique()
313
314
315 # In[159]:
316
317
318 unique, counts = np.unique(best_categ_shuffle, return_counts=True)
319 result = np.column_stack((unique, counts))
320 result
321
322
323 # In[160]:
324
325
326 df_t = pd.DataFrame()
327 df_t = pd.concat([df_t, df_input_x['cratio31']], axis = 1)
328 df_t = pd.concat([df_t, df_input_x['speed31']], axis = 1)
329 df_t = pd.concat([df_t, best_categ_shuffle], axis = 1)
330
331
332 # Balancear
333
334 # In[161]:
335
336
337 df_blc = pd.DataFrame()
338 for i in categorias:
339     df_i = df_t[df_t['best_categ'] == i]
340     df_i = df_i.head(n=9618)
341     df_blc = pd.concat([df_blc, df_i], axis = 0)
342
343
344 # In[162]:
345
346
347 dfinput_y = pd.DataFrame()
348 dfinput_y = pd.concat([dfinput_y, df_blc['best_categ']], axis = 1)
349
350
351 # In[163]:
352
353
354 dfinput_x = df_blc.drop(['best_categ'], axis = 1)
355
356

```

```

357 # Normalizar inputs
358
359 # In[164]:
360
361
362 dfcratios = pd.DataFrame()
363 dfcratios = pd.concat([dfcratios, dfinput_x.iloc[:, 0]], axis =0)
364
365
366 # In[165]:
367
368
369 desv_cratio = dfcratios.std()[0]
370 desv_cratio
371
372
373 # In[166]:
374
375
376 mean_cratio = dfcratios.mean()[0]
377 mean_cratio
378
379
380 # In[167]:
381
382
383 dfspeeds = pd.DataFrame()
384 dfspeeds = pd.concat([dfspeeds, dfinput_x.iloc[:, 1]], axis =0)
385
386
387 # In[168]:
388
389
390 mean_speed = dfspeeds.mean()[0]
391 mean_speed
392
393
394 # In[169]:
395
396
397 desv_speed = dfspeeds.std()[0]
398 desv_speed
399
400
401 # In[170]:
402
403
404 cratio = dfinput_x["cratio31"]
405 speed = dfinput_x["speed31"]
406 ncratio = (cratio - mean_cratio) / desv_cratio
407 nspeed = (speed - mean_speed) / desv_speed
408

```

```

409 dfinput_x_norm = pd.DataFrame()
410 dfinput_x_norm.insert(0, "ncratio", ncratio)
411 dfinput_x_norm.insert(1, "nspeed", nspeed)
412
413
414 # In[171]:
415
416
417 dfinput_x = dfinput_x_norm
418
419
420 # In[172]:
421
422
423 dfinput_x.shape
424
425
426 # In[173]:
427
428
429 dfinput_y.shape
430
431
432 # # RED NEURONAL (3 categor as)
433
434 # In[174]:
435
436
437 from sklearn.model_selection import train_test_split
438
439 # Split the data
440 x_train, x_test, y_train, y_test = train_test_split(dfinput_x, dfinput_y,
441     test_size=0.1, shuffle= True)
442
443 # In[175]:
444
445
446 x_train.shape
447
448
449 # In[176]:
450
451
452 x_test.shape
453
454
455 # In[177]:
456
457
458 y_test_nocod = y_test.copy()
459

```

```

460
461 # In[178]:
462
463
464 from keras.utils import to_categorical
465 y_train = to_categorical(y_train)
466 y_test = to_categorical(y_test)
467
468
469 # In[179]:
470
471
472 from keras import models, layers
473 import tensorflow
474
475 tensorflow.random.set_seed(100)
476
477 model = models.Sequential()
478 model.add(layers.Dense(2, activation='relu', input_shape=(2,)))
479 model.add(layers.Dense(2, activation='relu'))
480 model.add(layers.Dense(3, activation='softmax'))
481
482
483 # In[180]:
484
485
486 model.compile(optimizer='adam',
487               loss='categorical_crossentropy',
488               metrics=['acc'])
489
490
491 # In[181]:
492
493
494 history = model.fit(x_train, y_train, epochs = 60, validation_data = (x_test,
495                               y_test))
496
497 # In[182]:
498
499
500 _, ax = plt.subplots(1, 2, figsize=(10, 5), dpi=100)
501 ax[0].plot(history.history['acc'], 'r')
502 ax[0].plot(history.history['val_acc'], 'g')
503 ax[0].set_xlabel("Num of Epochs")
504 ax[0].set_ylabel("Accuracy")
505 ax[0].set_title("Training Accuracy vs Validation Accuracy")
506 ax[0].legend(['train', 'validation'])
507
508 ax[1].plot(history.history['loss'], 'r')
509 ax[1].plot(history.history['val_loss'], 'g')
510 ax[1].set_xlabel("Num of Epochs")

```

```

511 ax[1].set_ylabel("Loss")
512 ax[1].set_title("Training Loss vs Validation Loss")
513 ax[1].legend(['train', 'validation'])
514
515 plt.tight_layout()
516
517
518 # In[183]:
519
520
521 y_hat = model.predict(x=x_test)
522
523
524 # In[184]:
525
526
527 y_pred_label = np.argmax(y_hat, axis = 1)
528 y_true_label = np.argmax(y_test, axis = 1)
529
530
531 # In[185]:
532
533
534 unique, counts = np.unique(y_pred_label, return_counts=True)
535 result = np.column_stack((unique, counts))
536 result
537
538
539 # In[186]:
540
541
542 unique, counts = np.unique(y_true_label, return_counts=True)
543 result = np.column_stack((unique, counts))
544 result
545
546
547 # In[187]:
548
549
550 tabla = pd.crosstab(y_pred_label, y_true_label)
551 tabla
552
553
554 # # Analisis resultados
555
556 # In[188]:
557
558
559 df_correctos = pd.DataFrame(columns=['cratio', 'speed', 'category'])
560
561
562 # In[189]:

```

```

563
564
565 for i in range (0, len(y_pred_label)):
566     cratio_act = x_test.iloc[i, 0]
567     speed_act = x_test.iloc[i, 1]
568     categ_act = y_test_nocod.iloc[i, 0]
569
570     df_correctos = df_correctos.append({'cratio': cratio_act, 'speed':
571     speed_act, 'category': categ_act}, ignore_index=True)
572
573
574 # In[190]:
575
576
577 fig = px.scatter(df_correctos, "cratio", "speed", color = "category")
578 fig.show()
579
580
581 # GR FICAS POR CATEGOR A (Y SU ERROR)
582
583 # In[191]:
584
585
586 df_0 = pd.DataFrame(columns=['cratio', 'speed', 'category'])
587
588
589 # In[192]:
590
591
592 for i in range (0, len(y_pred_label)):
593     cratio_act = x_test.iloc[i, 0]
594     speed_act = x_test.iloc[i, 1]
595     categ_act = y_test_nocod.iloc[i, 0]
596     if y_pred_label[i] == y_true_label[i] and y_pred_label[i] == 0.0:
597         df_0 = df_0.append({'cratio': cratio_act, 'speed': speed_act, 'category
598         ': 0}, ignore_index=True)
599
600 for i in range (0, len(y_pred_label)):
601     cratio_act = x_test.iloc[i, 0]
602     speed_act = x_test.iloc[i, 1]
603     if y_pred_label[i] != y_true_label[i] and y_true_label[i] == 0.0:
604         df_0 = df_0.append({'cratio': cratio_act, 'speed': speed_act, 'category
605         ': -1}, ignore_index=True)
606
607
608 # In[193]:
609
610 fig = px.scatter(df_0, "cratio", "speed", color = "category")
611 fig.show()

```

```

612
613 # In[194]:
614
615
616 df_1 = pd.DataFrame(columns=['cratio', 'speed', 'category'])
617
618
619 # In[195]:
620
621
622 for i in range (0, len(y_pred_label)):
623     cratio_act = x_test.iloc[i, 0]
624     speed_act = x_test.iloc[i, 1]
625     categ_act = y_test_nocod.iloc[i, 0]
626     if y_pred_label[i] == y_true_label[i] and y_pred_label[i] == 1.0:
627         df_1 = df_1.append({'cratio': cratio_act, 'speed': speed_act, 'category
        ': 1}, ignore_index=True)
628
629 for i in range (0, len(y_pred_label)):
630     cratio_act = x_test.iloc[i, 0]
631     speed_act = x_test.iloc[i, 1]
632     if y_pred_label[i] != y_true_label[i] and y_true_label[i] == 1.0:
633         df_1 = df_1.append({'cratio': cratio_act, 'speed': speed_act, 'category
        ': -2}, ignore_index=True)
634
635
636 # In[196]:
637
638
639 fig = px.scatter(df_1, "cratio", "speed", color = "category")
640 fig.show()
641
642
643 # In[197]:
644
645
646 df_2 = pd.DataFrame(columns=['cratio', 'speed', 'category'])
647
648
649 # In[198]:
650
651
652 for i in range (0, len(y_pred_label)):
653     cratio_act = x_test.iloc[i, 0]
654     speed_act = x_test.iloc[i, 1]
655     categ_act = y_test_nocod.iloc[i, 0]
656     if y_pred_label[i] == y_true_label[i] and y_pred_label[i] == 2.0:
657         df_2 = df_2.append({'cratio': cratio_act, 'speed': speed_act, 'category
        ': 2}, ignore_index=True)
658
659 for i in range (0, len(y_pred_label)):
660     cratio_act = x_test.iloc[i, 0]

```



```

661     speed_act = x_test.iloc[i, 1]
662     if y_pred_label[i] != y_true_label[i] and y_true_label[i] == 2.0:
663         df_2 = df_2.append({'cratio': cratio_act, 'speed': speed_act, 'category
': -3}, ignore_index=True)
664
665
666 # In[199]:
667
668
669 fig = px.scatter(df_2, "cratio", "speed", color = "category")
670 fig.show()

```

B.2. Código Neural-Network-SPEED.ipynb

En este anexo se presenta el código realizado en Jupyter Notebook que prepara y manipula los datos, para finalmente construir y evaluar el modelo de red neuronal donde se optimiza el favour de max_speed. El fichero del código se llama **Neural-Network-SPEED.ipynb**.

```

1  #!/usr/bin/env python
2  # coding: utf-8
3
4  # In[120]:
5
6
7  import pandas as pd
8  import matplotlib.pyplot as plt
9  import plotly.express as px
10 import numpy as np
11
12
13 # In[121]:
14
15
16 from pathlib import Path
17 path = Path('../data_inputs_temp/')
18
19
20 # In[122]:
21
22
23 ficheros = ['entropy-shuffle', 'blosclz-shuffle-split', 'lz4-shuffle-split']
24
25 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
26 dataframes1 = {}
27 nsplits = [55296] * 3
28
29 for i in range(0, len(ficheros)):
30     dataframes1[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')

```

```

31     codec_filter = np.repeat(ficheros[i], nsplits[i])
32     dataframes1[ficheros[i]]['codec_filter'] = codec_filter
33     if i == 0:
34         categoria = np.repeat(str(30+i), nsplits[i])
35     else:
36         categoria = np.repeat(str(i-1), nsplits[i])
37     dataframes1[ficheros[i]]['categoria'] = categoria
38
39     dataframes1[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
40
41
42 # In[123]:
43
44
45 from pathlib import Path
46 path = Path('../data_inputs')
47
48
49 # In[124]:
50
51
52 ficheros = ['entropy-shuffle', 'blosclz-shuffle-split', 'lz4-shuffle-split']
53
54 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
55 dataframes2 = {}
56 nsplits = [55296] * 3
57
58 for i in range(0, len(ficheros)):
59     dataframes2[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
60     codec_filter = np.repeat(ficheros[i], nsplits[i])
61     dataframes2[ficheros[i]]['codec_filter'] = codec_filter
62     if i == 0:
63         categoria = np.repeat(str(30+i), nsplits[i])
64     else:
65         categoria = np.repeat(str(i-1), nsplits[i])
66     dataframes2[ficheros[i]]['categoria'] = categoria
67
68     dataframes2[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
69
70
71 # In[125]:
72
73
74 from pathlib import Path
75 path = Path('../data_inputs_wind/')
76
77
78 # In[126]:

```

```

79
80
81 ficheros = ['entropy-shuffle', 'blosclz-shuffle-split', 'lz4-shuffle-split']
82
83 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
84 dataframes3 = {}
85 nsplits = [55296] * 3
86
87 for i in range(0, len(ficheros)):
88     dataframes3[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
89     codec_filter = np.repeat(ficheros[i], nsplits[i])
90     dataframes3[ficheros[i]]['codec_filter'] = codec_filter
91     if i == 0:
92         categoria = np.repeat(str(30+i), nsplits[i])
93     else:
94         categoria = np.repeat(str(i-1), nsplits[i])
95     dataframes3[ficheros[i]]['categoria'] = categoria
96
97     dataframes3[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
98
99
100 # In[127]:
101
102
103 df_fig = pd.DataFrame()
104 for i in range(0, len(ficheros)):
105     df_fig = pd.concat([df_fig, dataframes1[ficheros[i]],
106                       dataframes2[ficheros[i]], dataframes3[ficheros[i]]],
107                       axis = 0)
108
109 # In[128]:
110
111
112 df_fig = df_fig[df_fig['special_vals'] == 0]
113
114
115 # In[129]:
116
117
118 fig = px.scatter(df_fig, "cratio", "speed", color = "categoria")
119 fig.show()
120
121
122 # 0: blosclz
123 #
124 # 1: lz4
125 #
126 # 30: entropy
127

```

```

128 # # Dataframe filter: SHUFFLE
129
130 # In[130]:
131
132
133 df_entropy = pd.DataFrame()
134 df_entropy = pd.concat([df_entropy, dataframes1[ficheros[0]], dataframes2[
135     ficheros[0]],
136     dataframes3[ficheros[0]]], axis = 0)
137 df_entropy.rename(columns = {'cratio':'cratio31', 'speed':'speed31', '
138     special_vals':'special_vals31'}, inplace = True)
139 df_entropy = df_entropy.drop(['nchunk', 'category', 'codec_filter', 'categoria
140     '], axis = 1)
141
142 # In[131]:
143
144 df_blosclz = pd.DataFrame()
145 df_blosclz = pd.concat([df_blosclz, dataframes1[ficheros[1]], dataframes2[
146     ficheros[1]],
147     dataframes3[ficheros[1]]], axis = 0)
148 df_blosclz.rename(columns = {'cratio':'cratio0', 'speed':'speed0', '
149     special_vals':'special_vals0'}, inplace = True)
150 df_blosclz = df_blosclz.drop(['nchunk', 'category', 'codec_filter', 'categoria
151     '], axis = 1)
152
153 # In[132]:
154
155 df_lz4 = pd.DataFrame()
156 df_lz4 = pd.concat([df_lz4, dataframes1[ficheros[2]], dataframes2[ficheros[2]],
157     dataframes3[ficheros[2]]], axis = 0)
158 df_lz4.rename(columns = {'cratio':'cratio1', 'speed':'speed1', 'special_vals':'
159     special_vals1'}, inplace = True)
160 df_lz4 = df_lz4.drop(['nchunk', 'category', 'codec_filter', 'categoria'], axis
161     = 1)
162
163 # In[133]:
164
165
166 df_shuffle = pd.DataFrame()
167 df_shuffle = pd.concat([df_shuffle, df_blosclz, df_lz4, df_entropy], axis = 1)
168
169
170 # Eliminar special_value (que lo sea en alguna de las 3 categor as)
171

```

```

172 # In[134]:
173
174
175 def delete_final(special_vals0, special_vals1, special_vals31):
176     if (special_vals0 == 1 or special_vals1 == 1 or special_vals31 == 1):
177         return 1
178     else:
179         return 0
180
181
182 # In[135]:
183
184
185 df_shuffle["delete_f"] = df_shuffle.apply(lambda x: delete_final(x[2], x[5], x
186     [8]), axis=1)
187
188 # In[136]:
189
190
191 df_shuffle = df_shuffle[df_shuffle['delete_f'] == 0]
192
193
194 # In[137]:
195
196
197 df_shuffle = df_shuffle.drop(['delete_f'], axis = 1)
198
199
200 # In[138]:
201
202
203 df_input_x = df_shuffle.iloc[:, [6, 7]]
204
205
206 # C lculo del SCORE
207
208 # In[139]:
209
210
211 df_shuffle = df_shuffle.drop(['cratio31', 'speed31', 'special_vals31', '
212     special_vals0', 'special_vals1'], axis = 1)
213
214 # In[140]:
215
216
217 dfcratios = pd.DataFrame()
218 for i in range(0, 2):
219     dfcratios = pd.concat([dfcratios, df_shuffle.iloc[:, 2*i]], axis =0)
220
221

```

```

222 # In[141]:
223
224
225 desv_cratio = dfcratios.std()[0]
226 desv_cratio
227
228
229 # In[142]:
230
231
232 mean_cratio = dfcratios.mean()[0]
233 mean_cratio
234
235
236 # In[143]:
237
238
239 dfspeeds = pd.DataFrame()
240 for i in range(0, 2):
241     dfspeeds = pd.concat([dfspeeds, df_shuffle.iloc[:, 2*i+1]], axis =0)
242
243
244 # In[144]:
245
246
247 mean_speed = dfspeeds.mean()[0]
248 mean_speed
249
250
251 # In[145]:
252
253
254 desv_speed = dfspeeds.std()[0]
255 desv_speed
256
257
258 # In[146]:
259
260
261 theta = 0.1
262
263 cratio0 = df_shuffle["cratio0"]
264 speed0 = df_shuffle["speed0"]
265 ncratio0 = (cratio0 - mean_cratio) / desv_cratio
266 nspeed0 = (speed0 - mean_speed) / desv_speed
267 df_shuffle.insert(2, 0, theta*ncratio0 + (1-theta)*nspeed0)
268
269 cratio1 = df_shuffle["cratio1"]
270 speed1 = df_shuffle["speed1"]
271 ncratio1 = (cratio1 - mean_cratio) / desv_cratio
272 nspeed1 = (speed1 - mean_speed) / desv_speed
273 df_shuffle.insert(5, 1, theta*ncratio1 + (1-theta)*nspeed1)

```

```

274
275
276 # In[147]:
277
278
279 df_pos = df_shuffle
280 df_pos = df_pos.drop(['cratio0', 'speed0', 'cratio1', 'speed1'], axis = 1)
281
282
283 # In[148]:
284
285
286 best_categ_shuffle = df_pos.idxmax(axis=1)
287 best_categ_shuffle = pd.DataFrame(best_categ_shuffle, columns=['best_categ'])
288
289
290 # In[149]:
291
292
293 categorias = best_categ_shuffle['best_categ'].unique()
294
295
296 # In[150]:
297
298
299 unique, counts = np.unique(best_categ_shuffle, return_counts=True)
300 result = np.column_stack((unique, counts))
301 result
302
303
304 # In[151]:
305
306
307 df_t = pd.DataFrame()
308 df_t = pd.concat([df_t, df_input_x['cratio31']], axis = 1)
309 df_t = pd.concat([df_t, df_input_x['speed31']], axis = 1)
310 df_t = pd.concat([df_t, best_categ_shuffle], axis = 1)
311
312
313 # Balancear
314
315 # In[152]:
316
317
318 df_blc = pd.DataFrame()
319 for i in categorias:
320     df_i = df_t[df_t['best_categ'] == i]
321     df_i = df_i.head(n=14154)
322     df_blc = pd.concat([df_blc, df_i], axis = 0)
323
324
325 # In[153]:

```

```

326
327
328 dfinput_y = pd.DataFrame()
329 dfinput_y = pd.concat([dfinput_y, df_blc['best_categ']], axis = 1)
330 dfinput_y.shape
331
332
333 # In[154]:
334
335
336 dfinput_x = df_blc.drop(['best_categ'], axis = 1)
337 dfinput_x.shape
338
339
340 # Normalizar inputs
341
342 # In[155]:
343
344
345 dfcratios = pd.DataFrame()
346 dfcratios = pd.concat([dfcratios, dfinput_x.iloc[:, 0]], axis =0)
347
348
349 # In[156]:
350
351
352 desv_cratio = dfcratios.std()[0]
353 desv_cratio
354
355
356 # In[157]:
357
358
359 mean_cratio = dfcratios.mean()[0]
360 mean_cratio
361
362
363 # In[158]:
364
365
366 dfspeeds = pd.DataFrame()
367 dfspeeds = pd.concat([dfspeeds, dfinput_x.iloc[:, 1]], axis =0)
368
369
370 # In[159]:
371
372
373 mean_speed = dfspeeds.mean()[0]
374 mean_speed
375
376
377 # In[160]:

```



```

378
379
380 desv_speed = dfspeeds.std()[0]
381 desv_speed
382
383
384 # In[161]:
385
386
387 cratio = dfinput_x["cratio31"]
388 speed = dfinput_x["speed31"]
389 ncratio = (cratio - mean_cratio) / desv_cratio
390 nspeed = (speed - mean_speed) / desv_speed
391
392 dfinput_x_norm = pd.DataFrame()
393 dfinput_x_norm.insert(0, "ncratio", ncratio)
394 dfinput_x_norm.insert(1, "nspeed", nspeed)
395
396
397 # In[162]:
398
399
400 dfinput_x = dfinput_x_norm
401
402
403 # In[163]:
404
405
406 dfinput_x.shape
407
408
409 # In[164]:
410
411
412 dfinput_y.shape
413
414
415 # # RED NEURONAL (2 categor as)
416
417 # In[165]:
418
419
420 from sklearn.model_selection import train_test_split
421
422 # Split the data
423 x_train, x_test, y_train, y_test = train_test_split(dfinput_x, dfinput_y,
424     test_size=0.1, shuffle= True)
425
426 # In[166]:
427
428

```

```

429 x_train.shape
430
431
432 # In[167]:
433
434
435 x_test.shape
436
437
438 # In[168]:
439
440
441 y_test_nocod = y_test.copy()
442
443
444 # In[169]:
445
446
447 from keras.utils import to_categorical
448 y_train = to_categorical(y_train)
449 y_test = to_categorical(y_test)
450
451
452 # In[170]:
453
454
455 from keras import models, layers
456 import tensorflow
457
458 tensorflow.random.set_seed(100)
459
460 model = models.Sequential()
461 model.add(layers.Dense(2, activation='relu', input_shape=(2,)))
462 model.add(layers.Dense(2, activation='relu'))
463 model.add(layers.Dense(2, activation='softmax'))
464
465
466 # In[171]:
467
468
469 model.compile(optimizer='adam',
470               loss='categorical_crossentropy',
471               metrics=['acc'])
472
473
474 # In[172]:
475
476
477 history = model.fit(x_train, y_train, epochs = 40, validation_data = (x_test,
478                               y_test))
479

```

```

480 # In[173]:
481
482
483 _, ax = plt.subplots(1, 2, figsize=(10, 5), dpi=100)
484 ax[0].plot(history.history['acc'], 'r')
485 ax[0].plot(history.history['val_acc'], 'g')
486 ax[0].set_xlabel("Num of Epochs")
487 ax[0].set_ylabel("Accuracy")
488 ax[0].set_title("Training Accuracy vs Validation Accuracy")
489 ax[0].legend(['train', 'validation'])
490
491 ax[1].plot(history.history['loss'], 'r')
492 ax[1].plot(history.history['val_loss'], 'g')
493 ax[1].set_xlabel("Num of Epochs")
494 ax[1].set_ylabel("Loss")
495 ax[1].set_title("Training Loss vs Validation Loss")
496 ax[1].legend(['train', 'validation'])
497
498 plt.tight_layout()
499
500
501 # In[174]:
502
503
504 y_hat = model.predict(x=x_test)
505
506
507 # In[175]:
508
509
510 y_pred_label = np.argmax(y_hat, axis = 1)
511 y_true_label = np.argmax(y_test, axis = 1)
512
513
514 # In[176]:
515
516
517 unique, counts = np.unique(y_pred_label, return_counts=True)
518 result = np.column_stack((unique, counts))
519 result
520
521
522 # In[177]:
523
524
525 unique, counts = np.unique(y_true_label, return_counts=True)
526 result = np.column_stack((unique, counts))
527 result
528
529
530 # In[178]:
531

```

```

532
533 tabla = pd.crosstab(y_pred_label, y_true_label)
534 tabla
535
536
537 # # Analisis resultados
538
539 # In[179]:
540
541
542 df_correctos = pd.DataFrame(columns=['cratio', 'speed', 'category'])
543
544
545 # In[180]:
546
547
548 for i in range (0, len(y_pred_label)):
549     cratio_act = x_test.iloc[i, 0]
550     speed_act = x_test.iloc[i, 1]
551     categ_act = y_test_nocod.iloc[i, 0]
552
553     df_correctos = df_correctos.append({'cratio': cratio_act, 'speed':
554     speed_act, 'category': categ_act}, ignore_index=True)
555
556
557 # In[181]:
558
559
560 fig = px.scatter(df_correctos, "cratio", "speed", color = "category")
561 fig.show()
562
563
564 # GR FICAS POR CATEGOR A (Y SU ERROR)
565
566 # In[182]:
567
568
569 df_0 = pd.DataFrame(columns=['cratio', 'speed', 'category'])
570
571
572 # In[183]:
573
574
575 for i in range (0, len(y_pred_label)):
576     cratio_act = x_test.iloc[i, 0]
577     speed_act = x_test.iloc[i, 1]
578     categ_act = y_test_nocod.iloc[i, 0]
579     if y_pred_label[i] == y_true_label[i] and y_pred_label[i] == 0.0:
580         df_0 = df_0.append({'cratio': cratio_act, 'speed': speed_act, 'category
581         ': 0}, ignore_index=True)

```

```

582 for i in range (0, len(y_pred_label)):
583     cratio_act = x_test.iloc[i, 0]
584     speed_act = x_test.iloc[i, 1]
585     if y_pred_label[i] != y_true_label[i] and y_true_label[i] == 0.0:
586         df_0 = df_0.append({'cratio': cratio_act, 'speed': speed_act, 'category
': -1}, ignore_index=True)
587
588
589 # In[184]:
590
591
592 fig = px.scatter(df_0, "cratio", "speed", color = "category")
593 fig.show()
594
595
596 # In[185]:
597
598
599 df_1 = pd.DataFrame(columns=['cratio', 'speed', 'category'])
600
601
602 # In[186]:
603
604
605 for i in range (0, len(y_pred_label)):
606     cratio_act = x_test.iloc[i, 0]
607     speed_act = x_test.iloc[i, 1]
608     categ_act = y_test_nocod.iloc[i, 0]
609     if y_pred_label[i] == y_true_label[i] and y_pred_label[i] == 1.0:
610         df_1 = df_1.append({'cratio': cratio_act, 'speed': speed_act, 'category
': 1}, ignore_index=True)
611
612 for i in range (0, len(y_pred_label)):
613     cratio_act = x_test.iloc[i, 0]
614     speed_act = x_test.iloc[i, 1]
615     if y_pred_label[i] != y_true_label[i] and y_true_label[i] == 1.0:
616         df_1 = df_1.append({'cratio': cratio_act, 'speed': speed_act, 'category
': -2}, ignore_index=True)
617
618
619 # In[187]:
620
621
622 fig = px.scatter(df_1, "cratio", "speed", color = "category")
623 fig.show()

```

B.3. Código Neural-Network-BALANCE.ipynb

En este anexo se presenta el código realizado en Jupyter Notebook que prepara y manipula los datos, para finalmente construir y evaluar el modelo de red neuronal donde se optimiza el favour de balance. El fichero del código se llama **Neural-Network-BALANCE.ipynb**.

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 # In[241]:
5
6
7 import pandas as pd
8 import matplotlib.pyplot as plt
9 import plotly.express as px
10 import numpy as np
11
12
13 # In[242]:
14
15
16 from pathlib import Path
17 path = Path('../data_inputs_temp/')
18
19
20 # In[243]:
21
22
23 ficheros = ['entropy-shuffle', 'blosclz-shuffle-split', 'zstd-shuffle-nosplit']
24
25 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
26 dataframes1 = {}
27 nsplits = [55296] * 2 + [13824]
28
29 for i in range(0, len(ficheros)):
30     dataframes1[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
31     codec_filter = np.repeat(ficheros[i], nsplits[i])
32     dataframes1[ficheros[i]]['codec_filter'] = codec_filter
33     if i == 0:
34         categoria = np.repeat(str(30+i), nsplits[i])
35     else:
36         categoria = np.repeat(str(i-1), nsplits[i])
37     dataframes1[ficheros[i]]['categoria'] = categoria
38
39     dataframes1[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
40
41
42 # In[244]:
43
```

```

44
45 from pathlib import Path
46 path = Path('../data_inputs')
47
48
49 # In[245]:
50
51
52 ficheros = ['entropy-shuffle', 'blosclz-shuffle-split', 'zstd-shuffle-nosplit']
53
54 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
55 dataframes2 = {}
56 nsplits = [55296] * 2 + [13824]
57
58 for i in range(0, len(ficheros)):
59     dataframes2[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
60     codec_filter = np.repeat(ficheros[i], nsplits[i])
61     dataframes2[ficheros[i]]['codec_filter'] = codec_filter
62     if i == 0:
63         categoria = np.repeat(str(30+i), nsplits[i])
64     else:
65         categoria = np.repeat(str(i-1), nsplits[i])
66     dataframes2[ficheros[i]]['categoria'] = categoria
67
68     dataframes2[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
69
70
71 # In[246]:
72
73
74 from pathlib import Path
75 path = Path('../data_inputs_wind/')
76
77
78 # In[247]:
79
80
81 ficheros = ['entropy-shuffle', 'blosclz-shuffle-split', 'zstd-shuffle-nosplit']
82
83 ficheros_path = [str(path) + '/' + fichero + '.csv.gz' for fichero in ficheros]
84 dataframes3 = {}
85 nsplits = [55296] * 2 + [13824]
86
87 for i in range(0, len(ficheros)):
88     dataframes3[ficheros[i]] = pd.read_csv(ficheros_path[i], delimiter=',')
89     codec_filter = np.repeat(ficheros[i], nsplits[i])
90     dataframes3[ficheros[i]]['codec_filter'] = codec_filter
91     if i == 0:
92         categoria = np.repeat(str(30+i), nsplits[i])
93     else:

```

```

94     categoria = np.repeat(str(i-1), nsplits[i])
95     dataframes3[ficheros[i]]['categoria'] = categoria
96
97     dataframes3[ficheros[i]].rename(columns = {'cratio':'cratio', ' speed':'
speed', ' special_vals':'special_vals', ' nchunk':'nchunk', ' category':'
category'}, inplace = True)
98
99
100 # In[248]:
101
102
103 55296/13824
104
105
106 # In[249]:
107
108
109 dataframes1[ficheros[2]] = dataframes1[ficheros[2]].loc[np.repeat(dataframes1[
ficheros[2]].index.values, 4)]
110 dataframes1[ficheros[2]].reset_index(drop=True)
111
112 dataframes2[ficheros[2]] = dataframes2[ficheros[2]].loc[np.repeat(dataframes2[
ficheros[2]].index.values, 4)]
113 dataframes2[ficheros[2]].reset_index(drop=True)
114
115 dataframes3[ficheros[2]] = dataframes3[ficheros[2]].loc[np.repeat(dataframes3[
ficheros[2]].index.values, 4)]
116 dataframes3[ficheros[2]].reset_index(drop=True)
117
118
119 # In[250]:
120
121
122 df_fig = pd.DataFrame()
123 for i in range(0, len(ficheros)):
124     df_fig = pd.concat([df_fig, dataframes1[ficheros[i]],
125                       dataframes2[ficheros[i]], dataframes3[ficheros[i]]],
126                       axis = 0)
127
128 # In[251]:
129
130
131 df_fig = df_fig[df_fig['special_vals'] == 0]
132
133
134 # In[252]:
135
136
137 fig = px.scatter(df_fig, "cratio", "speed", color = "categoria")
138 fig.show()
139

```



```

140
141 # 0: blosclz-split
142 #
143 # 1: zstd-nosplit
144 #
145 # 30: entropy
146
147 # # Dataframe filter: SHUFFLE
148
149 # In[253]:
150
151
152 df_entropy = pd.DataFrame()
153 df_entropy = pd.concat([df_entropy, dataframes1[ficheros[0]], dataframes2[
154     ficheros[0]],
155     dataframes3[ficheros[0]]], axis = 0)
156 df_entropy.rename(columns = {'cratio':'cratio31', 'speed':'speed31', '
157     special_vals':'special_vals31'}, inplace = True)
158 df_entropy = df_entropy.drop(['nchunk', 'category', 'codec_filter', 'categoria
159     '], axis = 1)
160
161 df_entropy = df_entropy.reset_index(drop=True)
162
163 # In[254]:
164
165 df_blosclz = pd.DataFrame()
166 df_blosclz = pd.concat([df_blosclz, dataframes1[ficheros[1]], dataframes2[
167     ficheros[1]],
168     dataframes3[ficheros[1]]], axis = 0)
169 df_blosclz.rename(columns = {'cratio':'cratio0', 'speed':'speed0', '
170     special_vals':'special_vals0'}, inplace = True)
171 df_blosclz = df_blosclz.drop(['nchunk', 'category', 'codec_filter', 'categoria
172     '], axis = 1)
173
174 df_blosclz = df_blosclz.reset_index(drop=True)
175
176 # In[255]:
177
178 df_zstd = pd.DataFrame()
179 df_zstd = pd.concat([df_zstd, dataframes1[ficheros[2]], dataframes2[ficheros
180     [2]],
181     dataframes3[ficheros[2]]], axis = 0)
182 df_zstd.rename(columns = {'cratio':'cratio1', 'speed':'speed1', 'special_vals
183     ':'special_vals1'}, inplace = True)

```

```

183 df_zstd = df_zstd.drop(['nchunk', 'category', 'codec_filter', 'categoria'],
184 axis = 1)
185 df_zstd = df_zstd.reset_index(drop=True)
186
187
188 # In[256]:
189
190
191 df_shuffle = pd.DataFrame()
192 df_shuffle = pd.concat([df_shuffle, df_blosclz, df_zstd, df_entropy], axis = 1)
193
194
195 # Eliminar special_value (que lo sea en alguna de las 3 categorias)
196
197 # In[257]:
198
199
200 def delete_final(special_vals0, special_vals1, special_vals31):
201     if (special_vals0 == 1 or special_vals1 == 1 or special_vals31 == 1):
202         return 1
203     else:
204         return 0
205
206
207 # In[258]:
208
209
210 df_shuffle["delete_f"] = df_shuffle.apply(lambda x: delete_final(x[2], x[5], x
211 [8]), axis=1)
212
213 # In[259]:
214
215
216 df_shuffle = df_shuffle[df_shuffle['delete_f'] == 0]
217
218
219 # In[260]:
220
221
222 df_shuffle = df_shuffle.drop(['delete_f'], axis = 1)
223
224
225 # In[261]:
226
227
228 df_input_x = df_shuffle.iloc[:, [6, 7]]
229
230
231 # C lculo del SCORE
232

```

```

233 # In[262]:
234
235
236 df_shuffle = df_shuffle.drop(['cratio31', 'speed31', 'special_vals31', '
    special_vals0', 'special_vals1'], axis = 1)
237
238
239 # In[263]:
240
241
242 dfcratios = pd.DataFrame()
243 for i in range(0, 2):
244     dfcratios = pd.concat([dfcratios, df_shuffle.iloc[:, 2*i]], axis =0)
245
246
247 # In[264]:
248
249
250 desv_cratio = dfcratios.std()[0]
251 desv_cratio
252
253
254 # In[265]:
255
256
257 mean_cratio = dfcratios.mean()[0]
258 mean_cratio
259
260
261 # In[266]:
262
263
264 dfspeeds = pd.DataFrame()
265 for i in range(0, 2):
266     dfspeeds = pd.concat([dfspeeds, df_shuffle.iloc[:, 2*i+1]], axis =0)
267
268
269 # In[267]:
270
271
272 mean_speed = dfspeeds.mean()[0]
273 mean_speed
274
275
276 # In[268]:
277
278
279 desv_speed = dfspeeds.std()[0]
280 desv_speed
281
282
283 # In[269]:

```

```

284
285
286 theta = 0.5
287
288 cratio0 = df_shuffle["cratio0"]
289 speed0 = df_shuffle["speed0"]
290 ncratio0 = (cratio0 - mean_cratio) / desv_cratio
291 nspeed0 = (speed0 - mean_speed) / desv_speed
292 df_shuffle.insert(2, 0, theta*ncratio0 + (1-theta)*nspeed0)
293
294 cratio1 = df_shuffle["cratio1"]
295 speed1 = df_shuffle["speed1"]
296 ncratio1 = (cratio1 - mean_cratio) / desv_cratio
297 nspeed1 = (speed1 - mean_speed) / desv_speed
298 df_shuffle.insert(5, 1, theta*ncratio1 + (1-theta)*nspeed1)
299
300
301 # In[270]:
302
303
304 df_pos = df_shuffle
305 df_pos = df_pos.drop(['cratio0', 'speed0', 'cratio1', 'speed1'], axis = 1)
306
307
308 # In[271]:
309
310
311 best_categ_shuffle = df_pos.idxmax(axis=1)
312 best_categ_shuffle = pd.DataFrame(best_categ_shuffle, columns=['best_categ'])
313
314
315 # In[272]:
316
317
318 categorias = best_categ_shuffle['best_categ'].unique()
319
320
321 # In[273]:
322
323
324 unique, counts = np.unique(best_categ_shuffle, return_counts=True)
325 result = np.column_stack((unique, counts))
326 result
327
328
329 # In[274]:
330
331
332 df_t = pd.DataFrame()
333 df_t = pd.concat([df_t, df_input_x['cratio31']], axis = 1)
334 df_t = pd.concat([df_t, df_input_x['speed31']], axis = 1)
335 df_t = pd.concat([df_t, best_categ_shuffle], axis = 1)

```

```

336
337
338 # Balancear
339
340 # In[275]:
341
342
343 df_blc = pd.DataFrame()
344 for i in categorias:
345     df_i = df_t[df_t['best_categ'] == i]
346     df_i = df_i.head(n=990)
347     df_blc = pd.concat([df_blc, df_i], axis = 0)
348
349
350 # In[276]:
351
352
353 dfinput_y = pd.DataFrame()
354 dfinput_y = pd.concat([dfinput_y, df_blc['best_categ']], axis = 1)
355 dfinput_y.shape
356
357
358 # In[277]:
359
360
361 dfinput_x = df_blc.drop(['best_categ'], axis = 1)
362 dfinput_x.shape
363
364
365 # Normalizar inputs
366
367 # In[278]:
368
369
370 dfcratios = pd.DataFrame()
371 dfcratios = pd.concat([dfcratios, dfinput_x.iloc[:, 0]], axis = 0)
372
373
374 # In[279]:
375
376
377 desv_cratio = dfcratios.std()[0]
378 desv_cratio
379
380
381 # In[280]:
382
383
384 mean_cratio = dfcratios.mean()[0]
385 mean_cratio
386
387

```

```

388 # In[281]:
389
390
391 dfspeeds = pd.DataFrame()
392 dfspeeds = pd.concat([dfspeeds, dfinput_x.iloc[:, 1]], axis =0)
393
394
395 # In[282]:
396
397
398 mean_speed = dfspeeds.mean()[0]
399 mean_speed
400
401
402 # In[283]:
403
404
405 desv_speed = dfspeeds.std()[0]
406 desv_speed
407
408
409 # In[284]:
410
411
412 cratio = dfinput_x["cratio31"]
413 speed = dfinput_x["speed31"]
414 ncratio = (cratio - mean_cratio) / desv_cratio
415 nspeed = (speed - mean_speed) / desv_speed
416
417 dfinput_x_norm = pd.DataFrame()
418 dfinput_x_norm.insert(0, "ncratio", ncratio)
419 dfinput_x_norm.insert(1, "nspeed", nspeed)
420
421
422 # In[285]:
423
424
425 dfinput_x = dfinput_x_norm
426
427
428 # In[286]:
429
430
431 dfinput_x.shape
432
433
434 # In[287]:
435
436
437 dfinput_y.shape
438
439

```

```

440 # # RED NEURONAL (2 categor as)
441
442 # In[288]:
443
444
445 from sklearn.model_selection import train_test_split
446
447 # Split the data
448 x_train, x_test, y_train, y_test = train_test_split(dfinput_x, dfinput_y,
449     test_size=0.1, shuffle= True)
450
451 # In[289]:
452
453
454 x_train.shape
455
456
457 # In[290]:
458
459
460 x_test.shape
461
462
463 # In[291]:
464
465
466 y_test_nocod = y_test.copy()
467
468
469 # In[292]:
470
471
472 from keras.utils import to_categorical
473 y_train = to_categorical(y_train)
474 y_test = to_categorical(y_test)
475
476
477 # In[293]:
478
479
480 from keras import models, layers
481 import tensorflow
482
483 tensorflow.random.set_seed(100)
484
485 model = models.Sequential()
486 model.add(layers.Dense(2, activation='relu', input_shape=(2,)))
487 model.add(layers.Dense(2, activation='relu'))
488 model.add(layers.Dense(2, activation='softmax'))
489
490

```

```

491 # In[294]:
492
493
494 model.compile(optimizer='adam',
495               loss='categorical_crossentropy',
496               metrics=['acc'])
497
498
499 # In[295]:
500
501
502 history = model.fit(x_train, y_train, epochs = 50, validation_data = (x_test,
503                             y_test))
504
505 # In[296]:
506
507
508 _, ax = plt.subplots(1, 2, figsize=(10, 5), dpi=100)
509 ax[0].plot(history.history['acc'], 'r')
510 ax[0].plot(history.history['val_acc'], 'g')
511 ax[0].set_xlabel("Num of Epochs")
512 ax[0].set_ylabel("Accuracy")
513 ax[0].set_title("Training Accuracy vs Validation Accuracy")
514 ax[0].legend(['train', 'validation'])
515
516 ax[1].plot(history.history['loss'], 'r')
517 ax[1].plot(history.history['val_loss'], 'g')
518 ax[1].set_xlabel("Num of Epochs")
519 ax[1].set_ylabel("Loss")
520 ax[1].set_title("Training Loss vs Validation Loss")
521 ax[1].legend(['train', 'validation'])
522
523 plt.tight_layout()
524
525
526 # In[297]:
527
528
529 y_hat = model.predict(x=x_test)
530
531
532 # In[298]:
533
534
535 y_pred_label = np.argmax(y_hat, axis = 1)
536 y_true_label = np.argmax(y_test, axis = 1)
537
538
539 # In[299]:
540
541

```



```

542 unique, counts = np.unique(y_pred_label, return_counts=True)
543 result = np.column_stack((unique, counts))
544 result
545
546
547 # In[300]:
548
549
550 unique, counts = np.unique(y_true_label, return_counts=True)
551 result = np.column_stack((unique, counts))
552 result
553
554
555 # In[301]:
556
557
558 tabla = pd.crosstab(y_pred_label, y_true_label)
559 tabla
560
561
562 # # Analisis resultados
563
564 # In[302]:
565
566
567 df_correctos = pd.DataFrame(columns=['cratio', 'speed', 'category'])
568
569
570 # In[303]:
571
572
573 for i in range(0, len(y_pred_label)):
574     cratio_act = x_test.iloc[i, 0]
575     speed_act = x_test.iloc[i, 1]
576     categ_act = y_test_nocod.iloc[i, 0]
577
578     df_correctos = df_correctos.append({'cratio': cratio_act, 'speed':
speed_act, 'category': categ_act}, ignore_index=True)
579
580
581
582 # In[304]:
583
584
585 fig = px.scatter(df_correctos, "cratio", "speed", color = "category")
586 fig.show()
587
588
589 # GR FICAS POR CATEGOR A (Y SU ERROR)
590
591 # In[305]:
592

```

```

593
594 df_0 = pd.DataFrame(columns=['cratio', 'speed', 'category'])
595
596
597 # In[306]:
598
599
600 for i in range (0, len(y_pred_label)):
601     cratio_act = x_test.iloc[i, 0]
602     speed_act = x_test.iloc[i, 1]
603     categ_act = y_test_nocod.iloc[i, 0]
604     if y_pred_label[i] == y_true_label[i] and y_pred_label[i] == 0.0:
605         df_0 = df_0.append({'cratio': cratio_act, 'speed': speed_act, 'category
': 0}, ignore_index=True)
606
607 for i in range (0, len(y_pred_label)):
608     cratio_act = x_test.iloc[i, 0]
609     speed_act = x_test.iloc[i, 1]
610     if y_pred_label[i] != y_true_label[i] and y_true_label[i] == 0.0:
611         df_0 = df_0.append({'cratio': cratio_act, 'speed': speed_act, 'category
': -1}, ignore_index=True)
612
613
614 # In[307]:
615
616
617 fig = px.scatter(df_0, "cratio", "speed", color = "category")
618 fig.show()
619
620
621 # In[308]:
622
623
624 df_1 = pd.DataFrame(columns=['cratio', 'speed', 'category'])
625
626
627 # In[309]:
628
629
630 for i in range (0, len(y_pred_label)):
631     cratio_act = x_test.iloc[i, 0]
632     speed_act = x_test.iloc[i, 1]
633     categ_act = y_test_nocod.iloc[i, 0]
634     if y_pred_label[i] == y_true_label[i] and y_pred_label[i] == 1.0:
635         df_1 = df_1.append({'cratio': cratio_act, 'speed': speed_act, 'category
': 1}, ignore_index=True)
636
637 for i in range (0, len(y_pred_label)):
638     cratio_act = x_test.iloc[i, 0]
639     speed_act = x_test.iloc[i, 1]
640     if y_pred_label[i] != y_true_label[i] and y_true_label[i] == 1.0:

```

```
641     df_1 = df_1.append({'cratio': cratio_act, 'speed': speed_act, 'category
        ': -2}, ignore_index=True)
642
643
644 # In[310]:
645
646
647 fig = px.scatter(df_1, "cratio", "speed", color = "category")
648 fig.show()
```