



TRABAJO FIN DE GRADO
INGENIERÍA INFORMÁTICA

Implementación optimizada sobre sistemas heterogéneos de algoritmos de Deep Learning para clasificación de imágenes

Autor

David Sánchez Pérez

Directores

José Miguel Mantas Ruiz



ESCUELA TÉCNICA SUPERIOR DE INGENIERÍAS INFORMÁTICA Y DE
TELECOMUNICACIÓN

Granada, mes de Febrero 2024



Implementación optimizada sobre
sistemas heterogéneos de algoritmos de
Deep Learning para clasificación de
imágenes

Autor

David Sánchez Pérez

Directores

José Miguel Mantas Ruiz

Título del Proyecto: Subtítulo del proyecto

Nombre Apellido1 Apellido2 (alumno)

Palabras clave: palabra_clave1, palabra_clave2, palabra_clave3,

Resumen

Poner aquí el resumen.

Project Title: Project Subtitle

First name, Family name (student)

Keywords: Keyword1, Keyword2, Keyword3,

Abstract

Write here the abstract in English.

Yo, **Nombre Apellido1 Apellido2**, alumno de la titulación TITULACIÓN de la **Escuela Técnica Superior de Ingenierías Informática y de Telecomunicación de la Universidad de Granada**, con DNI XXXXXXXXXX, autorizo la ubicación de la siguiente copia de mi Trabajo Fin de Grado en la biblioteca del centro para que pueda ser consultada por las personas que lo deseen.

Fdo: Nombre Apellido1 Apellido2

Granada a X de mes de 201 .

D. **Nombre Apellido1 Apellido2 (tutor1)**, Profesor del Área de XXXX del Departamento YYYY de la Universidad de Granada.

D. **Nombre Apellido1 Apellido2 (tutor2)**, Profesor del Área de XXXX del Departamento YYYY de la Universidad de Granada.

Informan:

Que el presente trabajo, titulado ***Título del proyecto, Subtítulo del proyecto***, ha sido realizado bajo su supervisión por **Nombre Apellido1 Apellido2 (alumno)**, y autorizamos la defensa de dicho trabajo ante el tribunal que corresponda.

Y para que conste, expiden y firman el presente informe en Granada a X de mes de 201 .

Los directores:

Nombre Apellido1 Apellido2 (tutor1) **Nombre Apellido1 Apellido2 (tutor2)**

Agradecimientos

Poner aquí agradecimientos...

Índice general

1. Introducción	1
1.1. Resumen	1
1.2. Estado del arte	1
1.3. Objetivos	1
2. Conceptos previos	3
2.1. Machine Learning	3
2.2. Componentes necesarios para el aprendizaje supervisado . . .	3
2.3. División de datos en entrenamiento y test	4
2.4. Tipos de aprendizaje	4
2.4.1. Aprendizaje Supervisado	4
2.4.2. Aprendizaje No Supervisado	4
2.4.3. Aprendizaje Por Refuerzo	4
2.5. Redes Neuronales Totalmente Conectadas	5
2.5.1. Neurona	5
2.5.2. Estructura por capas	5
2.5.3. Funciones de activación	6
2.5.4. One-hot encoding	8
2.5.5. Función de error o pérdida	8
2.5.6. Descenso del gradiente	8
2.5.7. Codificación One-Hot	9
2.5.8. Propagación hacia delante	9
2.5.9. Retropropagación	10
2.6. Redes Neuronales Convolucionales	12
2.6.1. Tipos de capas	12
2.6.2. Estructura por capas	12
2.6.3. ForwardPropagation	12
3. Aportaciones	13
3.1. Redes Neuronales Totalmente Conectadas	13
3.1.1. BackPropagation con 1 capa oculta	13
3.1.2. BackPropagation con 2 capas ocultas	18

4. Adaptación GPU	21
4.1. GPU en Redes Neuronales Totalmente Conectadas	21
4.2. GPU en Redes Neuronales Convolucionales	21
5. Comparación con distintas plataformas	23
5.1. cuDNN	23

Índice de figuras

2.1. Imagen de una neurona	5
2.2. Imagen de una capa de neuronas	5
2.3. Imagen de la función de activación ReLU	6
2.4. Imagen de la función de activación Sigmoides	7
2.5. Imagen de la función de activación SoftMax	7
2.6. Ejemplo de funcionamiento del descenso del gradiente	8
2.7. Estructura de una red con softmax	10
3.1. Red Neuronal totalmente conectada con 1 capa oculta	13
3.2. Imagen de backpropagation en la capa output	14
3.3. Imagen de backpropagation en los pesos entre la capa oculta y la capa output	16
3.4. Imagen de backpropagation en la capa oculta h1	16
3.5. Imagen de backpropagation en los pesos entre la capa input y la capa oculta	17
3.6. Imagen de backpropagation en la capa input	18
3.7. Red Neuronal totalmente conectada con 2 capas ocultas	18

Índice de cuadros

Capítulo 1

Introducción

1.1. Resumen

1.2. Estado del arte

1.3. Objetivos

Capítulo 2

Conceptos previos

2.1. Machine Learning

Se entiende como el campo de las ciencias de computación que en vez de enfocarse en el diseño de algoritmos explícitos, optan por el estudio de técnicas de aprendizaje. Este enfoque tiene un gran éxito en tareas computacionales donde no es factible diseñar un algoritmo de forma explícita. [1] En vez de averiguar las distintas reglas a seguir para llegar a una solución, esta alternativa permite simplemente suministrar ejemplos de lo que debería pasar en distintas situaciones, y dejar que la máquina aprenda y extraiga ella misma sus propias conclusiones. De esta forma, el procedimiento en aprendizaje supervisado consiste en 'entrenar' con una muestra de N ejemplos, extraer información de ellos, y posteriormente poder evaluar de forma 'correcta' (bajo un margen de error controlado) otra muestra de M ejemplos, siendo $M > N$. [2]

Este enfoque ha contribuido en el avance de áreas como reconocimiento de voz, visión por ordenador, procesamiento de lenguaje natural, etc.

2.2. Componentes necesarios para el aprendizaje supervisado

Datos de entrada X y de salida Y que el modelo empleará para aprender y tomar decisiones. Ambos se unen para formar un dataset de entradas-salidas $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$. Para que el aprendizaje sea posible, debe existir una función $F: X \rightarrow Y$ tal que $y_i = F(x_i)$ para $i \in \{1 \dots N\}$. De esta forma, en función del dataset D , el modelo tratará de encontrar una función G que aproxime F para dicho conjunto. Además, se suelen aplicar técnicas que permitan una mejor generalización del modelo, expandiendo las capacidades del mismo y permitiendo que su conocimiento pueda ser útil incluso fuera de la muestra de datos inicial. [2]

2.3. División de datos en entrenamiento y test

Para visualizar la generalización del modelo, el conjunto de datos D se suele dividir en 2 subconjuntos, (entrenamiento y test) de forma que se pueda estimar si realmente 'aprende' o solo memoriza.

Una vez realizada la división, se entrena el modelo con los datos del conjunto de entrenamiento. Cuando se termina el entrenamiento, se accede al conjunto test y se visualiza el rendimiento del modelo sobre el mismo. Como los datos de test no se emplearon en ningún momento, aportan una estimación sobre la generalización del modelo fuera de la muestra con la que se entrenó.

2.4. Tipos de aprendizaje

2.4.1. Aprendizaje Supervisado

Es el que se empleará en este proyecto. Se caracteriza por la presencia de una etiqueta 'correcta' y_i asociada a cada dato de entrada x_i . Posteriormente, la red empleará ambos valores para, a partir de x_i , tratar de deducir y_i . [2]

Aunque se trate de impedir, la existencia de ruido en los datos es inevitable, implicando que algunas etiquetas de $Y=\{y_1, y_2, \dots, y_N\}$ puedan ser erróneas. Este tipo de aprendizaje se divide a su vez en problemas de clasificación y regresión, centrándose en predecir etiquetas o valores numéricos, respectivamente.

2.4.2. Aprendizaje No Supervisado

En este tipo de aprendizaje los datos no contienen ninguna información respecto a Y . De esta forma, el conjunto de datos D se compone exclusivamente de valores $X=\{x_1, x_2, \dots, x_N\}$. [2]

2.4.3. Aprendizaje Por Refuerzo

En este caso tampoco existe un y_i 'correcto' asociado a cada x_i . En su lugar, se asocia a cada x_i una etiqueta con un valor posible de y_i , además de una medida que indica como de bueno es el mismo. [2]

2.5. Redes Neuronales Totalmente Conectadas

2.5.1. Neurona

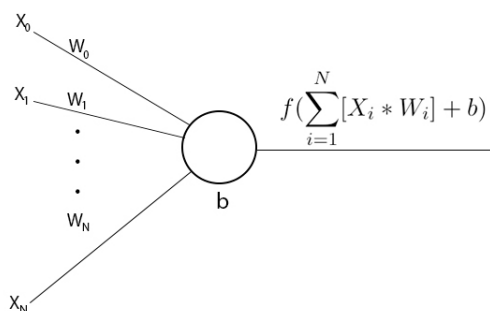


Figura 2.1: Imagen de una neurona

Una neurona parte de una serie de datos de entrada $X = \{x_1, x_2, \dots, x_N\}$ tal que cada $x_i \in X$ se encuentra asociado a un peso $w_i \in W$. Esta los emplea para realizar una suma ponderada y posteriormente añadir un sesgo b , además de aplicar una función de activación f sobre el resultado obtenido.

2.5.2. Estructura por capas

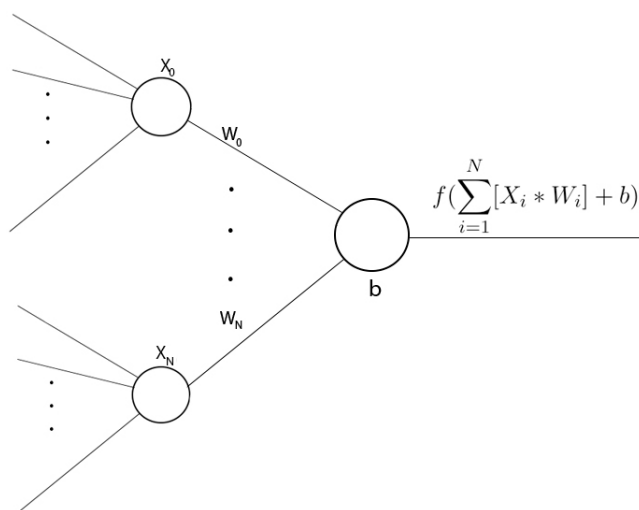


Figura 2.2: Imagen de una capa de neuronas

Las neuronas se suelen agrupar por capas, de tal forma que la salida de una compone la entrada de la siguiente, formando así modelos más sofisticados.

2.5.3. Funciones de activación

ReLU

$$ReLU(x) = \max(0, x) \quad (2.1)$$

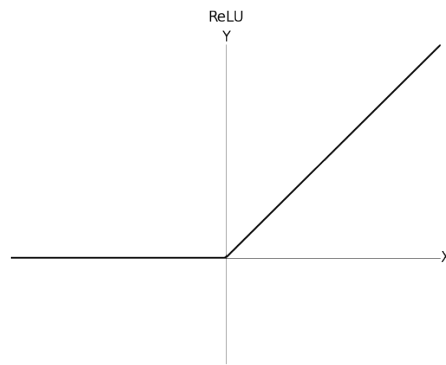


Figura 2.3: Imagen de la función de activación ReLU

A cambio de un bajo coste computacional, aporta no linealidad a la neurona, permitiendo a esta aprender funciones de mayor complejidad. Como su gradiente es 0 o 1, evita una reducción excesiva del mismo para valores positivos, mitigando así el problema del desvanecimiento del gradiente, caracterizado por la presencia de gradientes muy pequeños en backpropagation y provocar un aprendizaje lento. [3]

Sigmoide

$$sigmoide(x) = \frac{1}{1 + e^{-x}} \quad (2.2)$$



Figura 2.4: Imagen de la función de activación Sigmoide

Se trata de una función interesante en el ámbito de la clasificación binaria, pues se caracteriza por transformar un valor de entrada en una salida comprendida en el rango $[0-1]$.

Aunque sea monótona creciente y diferenciable en todos los puntos, tiende a saturarse con valores extremos (positivos o negativos). Por tanto, su aplicación dependerá del caso concreto a tratar. [4]

SoftMax

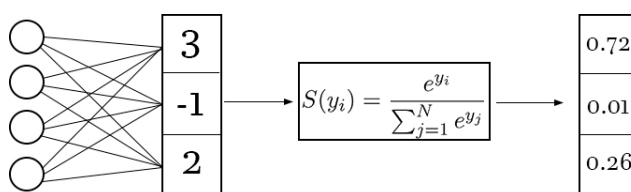


Figura 2.5: Imagen de la función de activación SoftMax

Para n entradas, produce n salidas con valores en el rango $[0-1]$ que mantienen la proporción de entrada y cuya suma es 1. Por tanto, se pueden interpretar como la probabilidad de pertenencia a cada clase, siendo especialmente útil en clasificación multiclase. [5]

2.5.4. One-hot encoding

2.5.5. Función de error o pérdida

Entropía Cruzada

$$E(y, \hat{y}) = - \sum_{i=1}^H [y_i * \log(\hat{y}_i)] \quad (2.3)$$

Es una métrica empleada en aprendizaje automático para medir qué tan bien se desempeña un modelo de clasificación. La pérdida o error se mide como un valor en el rango $[0-1]$, siendo 0 un modelo perfecto y 1 otro totalmente erróneo. [6]

H es el número de clases al que puede pertenecer cada dato de entrada $x_i \in X$.

Sigmoid Cross Entropy Loss

Es un caso particular de entropía cruzada caracterizado por la presencia de un número de clases igual a 2.

$$E(x) = -\frac{1}{N} \sum_{i=1}^N [y_i * \log(\hat{y}_i) + (1 - y_i) * \log(1 - \hat{y}_i)] \quad (2.4)$$

y = etiqueta real

\hat{y} = predicción

2.5.6. Descenso del gradiente

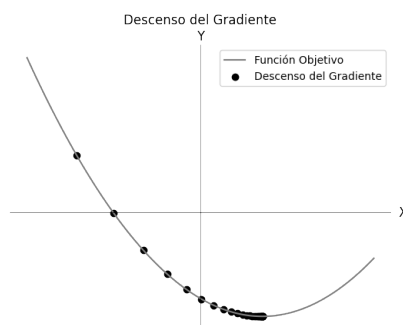


Figura 2.6: Ejemplo de funcionamiento del descenso del gradiente

Es un método de optimización iterativo que busca el mínimo local en una función diferenciable. En la figura 2.6 se muestra un ejemplo del mismo,

donde cada punto representa una iteración del algoritmo.

Su nombre viene del término 'gradiente', siendo este una generalización multivariable de la derivada y denominado por el símbolo ∇ . Para una función f y un punto p , este indica la dirección del máximo incremento en la misma. El descenso del gradiente usa esta información para, una vez obtenido el gradiente, desplazarse en dirección contraria, es decir, en dirección del mínimo. Además, la distancia que se recorre en cada iteración viene dada por un hiperparámetro denominado "learning rate" o α . [7] [8] [9]

Entrenamiento

De esta forma, el procedimiento para entrenar una red neuronal consiste en, para una entrada x_i y una etiqueta asociada y_i , emplear x_i para producir una predicción \hat{y}_i (**ForwardPropagation** o Propagación hacia delante) que posteriormente se podrá comparar con y_i mediante una función de error $H(x)$ y obtener una medida de lo buena o mala que fue la misma. Una vez obtenido dicho "error", se aplica el algoritmo del descenso del gradiente para cada parámetro de la red (**BackPropagation** o retropropagación). [6]

$$W_{t+1} = W_t - \alpha * \frac{\partial H(x)}{\partial W_t} \quad (2.5)$$

$$b_{t+1} = b_t - \alpha * \frac{\partial H(x)}{\partial b_t} \quad (2.6)$$

Así, se actualizarán los parámetros de la red neuronal según las fórmulas 2.5 y 2.6. En ellas, W_t y b_t indican los valores del peso W y bias b en el instante o iteración t , de la misma forma que W_{t+1} y b_{t+1} representan los valores de los mismos en el instante $t+1$. [10]

2.5.7. Codificación One-Hot

2.5.8. Propagación hacia delante

Se usará la codificación one-hot para cada $y_i \in Y$.

Se usará la función de error de entropía cruzada (Cross Entropy).

Se usará ReLU como función de activación en capas intermedias, y softmax para la capa output.

Suponemos que tenemos el caso $S(\hat{y}) = [S(\hat{y}_1), S(\hat{y}_2), S(\hat{y}_3)] = [0.04, 0.7, 0.26]$

Suponemos que para ese x_i caso dado, tenemos $y_i = [0, 0, 1]$

La red piensa que x_i pertenece a la clase 2 (0.7 es el mayor número del vector $S(\hat{y})$). Sin embargo, y_i dice que x_i pertenece a la clase 3.

Calculamos el error de entropía cruzada

$$E(y, S(\hat{y})) = -(0 * \log(0.04) + 0 * \log(0.7) + 1 * \log(0.26)) \quad (2.7)$$

$$E(y, S(\hat{y})) = -\log(0.26) = 0.585 \quad (2.8)$$

[11]

2.5.9. Retropopagación

Denotaremos como O_i a las salidas de la capa softmax. Las entradas a la capa softmax se denotarán como Z_i .

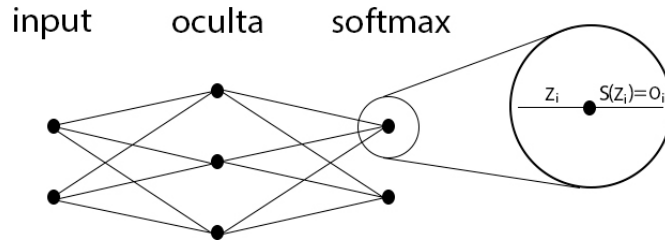


Figura 2.7: Estructura de una red con softmax

$$E = - \sum_{i=1}^C [y_i * \log(O_i)] \quad (2.9)$$

$$\frac{\partial E}{\partial Z_k} = \frac{\partial (- \sum_{i=1}^C [y_i * \log(O_i)])}{\partial Z_k} \quad (2.10)$$

$$\frac{\partial E}{\partial Z_k} = - \sum_{i=1}^C [\frac{\partial (y_i * \log(O_i))}{\partial Z_k}] \quad (2.11)$$

$$\frac{\partial E}{\partial Z_k} = - \sum_{i=1}^C [y_i * \frac{\partial (\log(O_i))}{\partial Z_k}] \quad (2.12)$$

$$\frac{\partial E}{\partial Z_k} = - \sum_{i=1}^C [y_i * \frac{\partial (\log(O_i))}{\partial O_i} * \frac{\partial O_i}{\partial Z_k}] \quad (2.13)$$

$$\frac{\partial E}{\partial Z_k} = - \sum_{i=1}^C [\frac{y_i}{O_i} * \frac{\partial O_i}{\partial Z_k}] \quad (2.14)$$

$$(2.15)$$

$$\frac{\partial h(x)}{\partial x} = \frac{f'(x) * g(x) - g'(x) * f(x)}{g(x)^2} \quad (2.16)$$

$$S(z_i) = \frac{e^{Z_i}}{e^{Z_1} + \dots + e^{Z_C}} \quad (2.17)$$

$$(2.18)$$

Hay 2 casos posibles, $\frac{S(Z_i)}{Z_i}$ o $\frac{S(Z_i)}{Z_j}$, donde $i \neq j$.

Caso $\frac{S(Z_i)}{Z_i}$

$$\frac{\partial S(Z_1)}{\partial Z_1} = \frac{[\frac{\partial e^{Z_1}}{\partial Z_1} * (e^{Z_1} + \dots + e^{Z_C})] - [\frac{\partial(e^{Z_1} + \dots + e^{Z_C})}{\partial Z_1} * e^{Z_1}]}{(e^{Z_1} + \dots + e^{Z_C})^2} \quad (2.19)$$

$$\frac{\partial S(Z_1)}{\partial Z_1} = \frac{[e^{Z_1} * \sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}]] - [e^{Z_1} * e^{Z_1}]}{(\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}])^2} \quad (2.20)$$

$$\frac{\partial S(Z_1)}{\partial Z_1} = \frac{e^{Z_1}(\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}] - e^{Z_1})}{(\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}])^2} \quad (2.21)$$

$$\frac{\partial S(Z_1)}{\partial Z_1} = \frac{e^{Z_1}}{\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}]} * \frac{\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}] - e^{Z_1}}{\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}]} \quad (2.22)$$

$$\frac{\partial S(Z_1)}{\partial Z_1} = S(Z_1) * (1 - S(Z_1)) \quad (2.23)$$

Caso $\frac{S(Z_i)}{Z_j}$

$$\frac{\partial S(Z_2)}{\partial Z_1} = \frac{[\frac{\partial e^{Z_2}}{\partial Z_1} * (e^{Z_1} + \dots + e^{Z_C})] - [\frac{\partial(e^{Z_1} + \dots + e^{Z_C})}{\partial Z_1} * e^{Z_2}]}{(e^{Z_1} + \dots + e^{Z_C})^2} \quad (2.24)$$

$$\frac{\partial S(Z_2)}{\partial Z_1} = \frac{[0 * \sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}]] - [e^{Z_1} * e^{Z_2}]}{(\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}])^2} \quad (2.25)$$

$$\frac{\partial S(Z_1)}{\partial Z_1} = \frac{-e^{Z_1} * e^{Z_2}}{(\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}])^2} \quad (2.26)$$

$$\frac{\partial S(Z_2)}{\partial Z_1} = \frac{-e^{Z_1}}{\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}]} * \frac{e^{Z_2}}{\sum_{i=1}^C [e^{Z_1} + \dots + e^{Z_C}]} \quad (2.27)$$

$$\frac{\partial S(Z_2)}{\partial Z_1} = -S(Z_1) * S(Z_2) \quad (2.28)$$

Ahora juntamos los 2 posibles casos.
 Parte izquierda cuando $i \neq k$, parte derecha cuando $i = k$.

$$\frac{\partial E}{\partial Z_k} = -\left[\sum_{i \neq k}^H \left[\frac{y_i}{O_i} * -O_i * O_k\right] + \frac{y_k}{O_k} * O_k * (1 - O_k)\right] \quad (2.29)$$

Simplificamos.

$$\frac{\partial E}{\partial Z_k} = -\left[\sum_{i \neq k}^H [-y_i * O_k] + [y_k * (1 - O_k)]\right] \quad (2.30)$$

$$\frac{\partial E}{\partial Z_k} = -\left[-O_k \sum_{i \neq k}^H -y_i + [y_k * (1 - O_k)]\right] \quad (2.31)$$

Como usamos codificación one-hot en Y.

$$\sum_{i=1}^H y_i = 1 \quad (2.32)$$

$$\sum_{i \neq k}^H y_i = \sum_{i=1}^H y_i - y_k = 1 - y_k \quad (2.33)$$

Usamos esto para simplificar lo que teníamos.

$$\frac{\partial E}{\partial Z_k} = -\left[-O_k \sum_{i \neq k}^H -y_i + [y_k * (1 - O_k)]\right] \quad (2.34)$$

$$\frac{\partial E}{\partial Z_k} = [O_k * (1 - y_k)] - [y_k * (1 - O_k)] \quad (2.35)$$

$$\frac{\partial E}{\partial Z_k} = O_k - O_k * y_k - y_k + O_k * y_k \quad (2.36)$$

$$\frac{\partial E}{\partial Z_k} = O_k - y_k \quad (2.37)$$

2.6. Redes Neuronales Convolucionales

2.6.1. Tipos de capas

2.6.2. Estructura por capas

2.6.3. ForwardPropagation

Capítulo 3

Aportaciones

3.1. Redes Neuronales Totalmente Conectadas

3.1.1. BackPropagation con 1 capa oculta

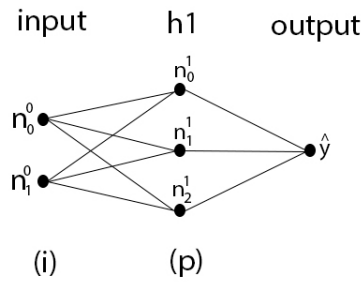


Figura 3.1: Red Neuronal totalmente conectada con 1 capa oculta

La Figura 3.1 se compone de puntos y líneas, representando neuronas y pesos que las conectan respectivamente. Cada punto es una neurona, y cada línea un peso.

La Figura 3.1 presenta 3 capas (input, h1, output) que corresponden a capa de entrada, capa oculta h_1 , y capa de salida respectivamente. El superíndice indica la capa a la que pertenece una neurona o peso, mientras que el subíndice indica el número del mismo en su respectiva capa. En el caso de los pesos, se requieren 2 subíndices para identificar a cada uno (pues un peso une 2 neuronas).

La capa de entrada se compone de 2 neuronas (n_0^0 y n_1^0).

La capa oculta h_1 tiene 3 neuronas (n_0^1 , n_1^1 , y n_2^1)

El peso W_{jk}^i referencia al peso que une las neuronas n_j^i y n_k^{i+1} .

Además, se denotará como \hat{y} a la última neurona de la red, pues contendrá la predicción de la misma.

Capa output

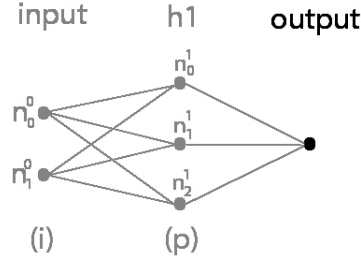


Figura 3.2: Imagen de backpropagation en la capa output

Sea la neurona n_j^i , se define como a_j^i el valor de dicha neurona antes de aplicar sobre ella su función de activación asociada, y z_j^i el obtenido tras aplicarla. De forma adicional, se usará \hat{a} y \hat{z} para \hat{y} .

Así, la función de pérdida 2.4 se convierte en:

$$H(x) = -\frac{1}{N} \sum_{i=1}^N [y_i * \log(\hat{z}_i) + (1 - y_i) * \log(1 - \hat{z}_i)] \quad (3.1)$$

Para realizar el descenso del gradiente, se debe empezar calculando la derivada de la función de pérdida respecto a la predicción obtenida. Es decir, la derivada de la fórmula (3.1) respecto de las neuronas en la última capa de la red tras aplicar sus respectivas funciones de activación, que en este caso corresponde a \hat{z} .

Por simplicidad, podemos dividir esta derivada en 2 partes.

Parte izquierda:

$$f(x) = A * B \quad (3.2)$$

$$f'(x) = AB' + A'B \quad (3.3)$$

$$\frac{\partial y_i}{\partial \hat{z}_i} = 0 \quad (3.4)$$

$$\frac{\partial \log(\hat{z}_i)}{\partial \hat{z}_i} = \frac{1}{\hat{z}_i} \quad (3.5)$$

$$\frac{\partial y_i * \log(\hat{z}_i)}{\partial \hat{z}_i} = y_i * \frac{1}{\hat{z}_i} + 0 * \log(\hat{z}_i) = \frac{y_i}{\hat{z}_i} \quad (3.6)$$

Parte derecha:

$$\frac{\partial(1 - y_i)}{\partial \hat{z}_i} = 0 \quad (3.7)$$

$$\frac{\partial \log(1 - \hat{z}_i)}{\partial \hat{z}_i} = \frac{1}{1 - \hat{z}_i} * (-1) \quad (3.8)$$

$$\frac{\partial(1 - y_i) * \log(1 - \hat{z}_i)}{\partial \hat{z}_i} = (1 - y_i) * \frac{1}{1 - \hat{z}_i} * (-1) + 0 * \log(1 - \hat{z}_i) \quad (3.9)$$

$$\frac{\partial(1 - y_i) * \log(1 - \hat{z}_i)}{\partial \hat{z}_i} = -\frac{1 - y_i}{1 - \hat{z}_i} \quad (3.10)$$

Finalmente, se obtiene:

$$\frac{\partial H(x)}{\partial \hat{z}_i} = -\frac{1}{N} \sum_{i=1}^N \left[\frac{y_i}{\hat{z}_i} - \frac{1 - y_i}{1 - \hat{z}_i} \right] \quad (3.11)$$

Función activación de la capa output

En la capa output se emplea la función de activación sigmoide.

$$\text{sigmoide}(x) = \frac{1}{1 + e^{-x}} \quad (3.12)$$

$$\text{sigmoide}'(x) = \frac{\text{sigmoide}(x)}{1 - \text{sigmoide}(x)} \quad (3.13)$$

De esta forma,

$$\frac{\partial \hat{z}}{\partial \hat{a}} = \frac{\partial \text{sigmoide}(\hat{a})}{\partial \hat{a}} = \text{sigmoide}(\hat{a}) * (1 - \text{sigmoide}(\hat{a})) \quad (3.14)$$

Ahora, podemos calcular el gradiente completo hasta la capa output antes de aplicar su función de activación.

$$\text{grad_output} = \frac{\partial H(x)}{\partial \hat{z}} * \frac{\partial \hat{z}}{\partial \hat{a}} = -\frac{1}{N} \sum_{i=1}^N \left[\frac{y_i}{\hat{z}_i} - \frac{1 - y_i}{1 - \hat{z}_i} \right] * \frac{\partial \hat{z}}{\partial \hat{a}} \quad (3.15)$$

$$\text{grad_output} = -\frac{1}{N} \sum_{i=1}^N \left[\frac{y_i}{\hat{z}_i} - \frac{1 - y_i}{1 - \hat{z}_i} \right] * \text{sigmoide}(\hat{a}) * (1 - \text{sigmoide}(\hat{a})) \quad (3.16)$$

Pesos capas h1-output

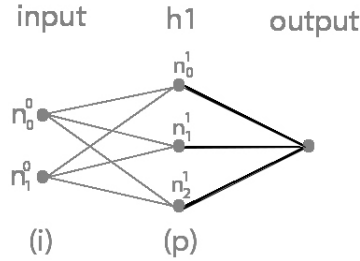


Figura 3.3: Imagen de backpropagation en los pesos entre la capa oculta y la capa output

Una vez calculado el gradiente hasta la capa output, se puede calcular el gradiente respecto a cada peso que se encuentra conectado a esta desde la capa anterior. Es decir, para cada $h_p^1 \in h_1$, se calcula $\frac{dH(x)}{dW_{p\hat{y}}^1}$. Usando la regla de la cadena, equivale a realizar lo siguiente:

$$\frac{\partial \hat{y}}{\partial W_{p\hat{y}}^1} = \frac{\partial n_p^1 * W_{p\hat{y}}^1}{\partial W_{p\hat{y}}^1} = n_p^1 \quad (3.17)$$

$$\frac{\partial H(x)}{\partial W_{p\hat{y}}^1} = \frac{\partial H(x)}{\partial \hat{y}} * \frac{\partial \hat{z}}{\partial \hat{a}} * \frac{\partial \hat{y}}{\partial W_{p\hat{y}}^1} = grad_output * \frac{\partial \hat{y}}{\partial W_{p\hat{y}}^1} = grad_output * n_p^1 \quad (3.18)$$

Capa oculta h1

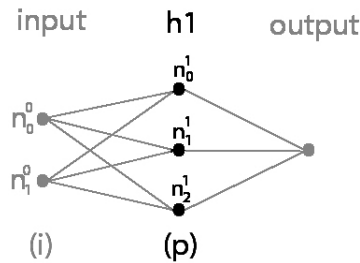


Figura 3.4: Imagen de backpropagation en la capa oculta h1

$$\frac{\partial \hat{y}}{\partial n_p^1} = \frac{\partial n_p^1 * \partial W_{p\hat{y}}^1}{\partial n_p^1} = \partial W_{p\hat{y}}^1 \quad (3.19)$$

$$\text{deriv_relu}(x) = 0 \text{ si } x \leq 0, 1 \text{ en caso contrario} \quad (3.20)$$

$$\frac{\partial z_p^1}{\partial a_p^1} = \frac{\partial \text{relu}(a_p^1)}{\partial a_p^1} = \text{deriv_relu}(a_p^1) \quad (3.21)$$

$$\frac{\partial H(x)}{\partial n_p^1} = \text{grad_output} * \frac{\partial \hat{y}}{\partial n_p^1} * \frac{\partial z_p^1}{\partial a_p^1} = \text{grad_output} * W_{p\hat{y}}^1 * \text{deriv_relu}(a_p^1) \quad (3.22)$$

$$\frac{\partial H(x)}{\partial n_p^1} = \text{gradient_n}_p^1 \quad (3.23)$$

Pesos capas input-h1

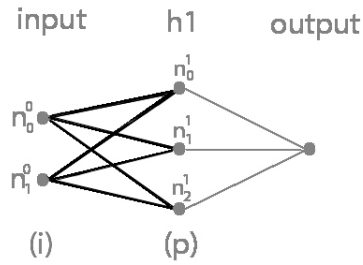


Figura 3.5: Imagen de backpropagation en los pesos entre la capa input y la capa oculta

$$\frac{\partial n_p^1}{\partial W_{ip}^0} = \frac{\partial n_i^0 * W_{ip}^0}{\partial W_{ip}^0} = n_i^0 \quad (3.24)$$

$$\frac{\partial H(x)}{\partial W_{ip}^0} = \text{gradient_n}_p^1 * \frac{\partial n_p^1}{\partial W_{ip}^0} = \text{gradient_n}_p^1 * n_i^0 \quad (3.25)$$

Capa input

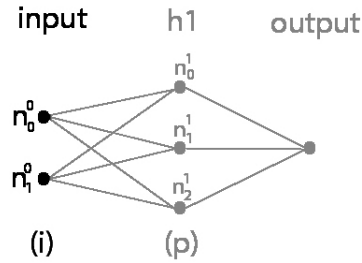


Figura 3.6: Imagen de backpropagation en la capa input

$$\frac{\partial n_p^1}{\partial n_i^0} = \frac{\partial n_i^0 * W_{ip}^0}{\partial n_i^0} = W_{ip}^0 \quad (3.26)$$

$$\frac{\partial z_i^0}{\partial a_i^0} = 1 \quad (3.27)$$

Ahora se realiza una sumatoria con todos los 'caminos posibles' hacia la misma neurona

$$\frac{\partial H(x)}{\partial n_i^0} = \sum_{p=0}^{h_1.size()} gradient_n_p^1 * \frac{\partial n_p^1}{\partial n_i^0} * \frac{\partial z_i^0}{\partial a_i^0} \quad (3.28)$$

$$\frac{\partial H(x)}{\partial n_i^0} = \sum_{p=0}^{h_1.size()} gradient_n_p^1 * W_{ip}^0 \quad (3.29)$$

3.1.2. BackPropagation con 2 capas ocultas

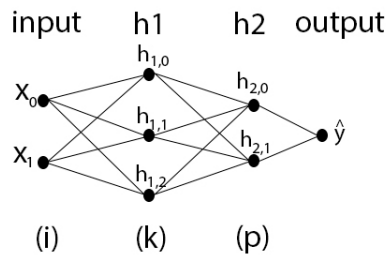


Figura 3.7: Red Neuronal totalmente conectada con 2 capas ocultas

Capa output

Se parte de la misma función de pérdida que en el caso anterior.

$$H(x) = -\frac{1}{N} \sum_{i=1}^N [y_i * \log(\hat{y}_{out_i}) + (1 - y_i) * \log(1 - \hat{y}_{out_i})] \quad (3.30)$$

$$\frac{dH(x)}{d\hat{y}_{out_i}} = -\frac{1}{N} \sum_{i=1}^N \left[\frac{y_i}{\hat{y}_{out_i}} - \frac{1 - y_i}{1 - \hat{y}_{out_i}} \right] \quad (3.31)$$

Función activación de la capa output

$$\frac{d\hat{y}_{out}}{d\hat{y}_{in}} = \frac{dsigmoide(\hat{y}_{in})}{d\hat{y}_{in}} = sigmoide(\hat{y}_{in}) * (1 - sigmoide(\hat{y}_{in})) \quad (3.32)$$

$$grad_output = -\frac{1}{N} \sum_{i=1}^N \left[\frac{y_i}{\hat{y}_{out_i}} - \frac{1 - y_i}{1 - \hat{y}_{out_i}} \right] * sigmoide(\hat{y}_{in}) * (1 - sigmoide(\hat{y}_{in})) \quad (3.33)$$

Pesos capas h1-h2

$$\frac{d\hat{y}_{in}}{dW_{p\hat{y}_{in}}} = \frac{dh_{2out_p} * W_{p\hat{y}_{in}}}{dW_{p\hat{y}_{in}}} = h_{2out_p} \quad (3.34)$$

$$\frac{dH(x)}{dW_{p\hat{y}_{in}}} = \frac{dH(x)}{d\hat{y}_{out}} * \frac{d\hat{y}_{out}}{d\hat{y}_{in}} * \frac{d\hat{y}_{in}}{dW_{p\hat{y}_{in}}} = grad_output * \frac{d\hat{y}_{in}}{dW_{p\hat{y}_{in}}} = grad_output * h_{2out_p} \quad (3.35)$$

Capa oculta h2

$$\frac{d\hat{y}_{in}}{dh_{2out_p}} = \frac{dh_{2out_p} * dW_{p\hat{y}_{in}}}{dh_{2out_p}} = dW_{p\hat{y}_{in}} \quad (3.36)$$

$$\frac{dh_{2out_p}}{dh_{2in_p}} = \frac{drelu(h_{2in_p})}{dh_{2in_p}} = 0 \text{ si } h_{2in_p} \leq 0, 1 \text{ en caso contrario} = deriv_relu(h_{2in_p}) \quad (3.37)$$

$$\frac{dH(x)}{dh_{2in_p}} = grad_output * \frac{d\hat{y}_{in}}{dh_{2out_p}} * \frac{dh_{2out_p}}{dh_{2in_p}} = grad_output * W_{p\hat{y}_{in}} * deriv_relu(h_{2in_p}) \quad (3.38)$$

$$\frac{dH(x)}{dh_{2in_p}} = grad_h2p_output \quad (3.39)$$

Pesos capas h1-h2

$$\frac{h_{2in_p}}{W_{kp}} = \frac{dh_{1out_k} * W_{kp}}{dW_{kp}} = h_{1out_k} \quad (3.40)$$

$$\frac{dH(x)}{dW_{kp}} = grad_h2p_output * \frac{dh_{2in_p}}{dW_{kp}} = grad_h2p_output * h_{1out_k} \quad (3.41)$$

Capa oculta h1

$$\frac{dh_{2in_p}}{dh_{1out_k}} = \frac{dh_{1out_k} * W_{kp}}{dh_{1out_k}} = W_{kp} \quad (3.42)$$

$$\frac{dh_{1out_k}}{dh_{1in_k}} = 1 \quad (3.43)$$

Ahora se realiza una sumatoria con todos los 'caminos posibles' hacia la misma neurona

$$\frac{dH(x)}{dh_{1in_k}} = \sum_{p=0}^{h2.size()} grad_h2p_output * \frac{dh_{2in_p}}{dh_{1out_k}} * \frac{dh_{1out_k}}{dh_{1in_k}} \quad (3.44)$$

$$\frac{dH(x)}{dh_{1in_k}} = deriv_relu(h_{1in_k}) * \sum_{p=0}^{h2.size()} grad_h2p_output * W_{kp} \quad (3.45)$$

$$grad_h1k_output = deriv_relu(h_{1in_k}) * \sum_{p=0}^{h2.size()} grad_h2p_output * W_{kp} \quad (3.46)$$

Capítulo 4

Adaptación GPU

- 4.1. GPU en Redes Neuronales Totalmente Conectadas
- 4.2. GPU en Redes Neuronales Convolucionales

Capítulo 5

Comparación con distintas plataformas

5.1. cuDNN

Bibliografía

- [1] Izzat El Hajj Wen-emi W.Hwu, David B.kirk. *Programming Massively Parallel Processors*. Morgan Kaufmann, 50 Hampshire Street, 5th Floor, Cambridge, MA 02139, United States, 4 edition, 2022.
- [2] Hsuan-Tien Lin Yaser S. Abu-Mostafa, Malik Magdon-Ismail. *Learning From Data*. California Institute of Technology Pasadena, CA 91125, USA, 1 edition, 2012.
- [3] Douglas Karr. Unidad lineal rectificada, 2024. <https://es.martech.zone/acronym/relu/> [Accessed:25/02/2024].
- [4] Javi. La función sigmoide: Una herramienta clave en redes neuronales, 2023. <https://jacar.es/la-funcion-sigmoide-una-herramienta-clave-en-redes-neuronales/> [Accessed:25/02/2024].
- [5] Jason Brownlee. Softmax activation function with python, 2020. <https://machinelearningmastery.com/softmax-activation-function-with-python/> [Accessed:24/02/2024].
- [6] Saurav Maheshkar. What is cross entropy loss? a tutorial with code, 2023. <https://wandb.ai/sauravmaheshkar/cross-entropy/reports/What-Is-Cross-Entropy-Loss-A-Tutorial-With-Code--VmlldzoxMDA5NTMx#:~:text=Cross%20entropy%20loss%20is%20a,close%20to%200%20as%20possible.> [Accessed:29/02/2024].
- [7] Descenso del gradiente, 2024. https://es.wikipedia.org/wiki/Descenso_del_gradiente#:~:text=El%20descenso%20del%20gradiente%20o,en%20direcci%C3%B3n%20contraria%20al%20gradiente. [Accessed:26/02/2024].
- [8] Gradiente, 2024. <https://es.wikipedia.org/wiki/Gradiente> [Accessed:26/02/2024].

-
- [9] Robert Kwiatkowski. Gradient descent algorithm — a deep dive, 2021. [https://towardsdatascience.com/gradient-descent-algorithm-a-deep-dive-cf04e8115f21#:~:text=Gradient%20descent%20\(GD\)%20is%20an,e.g.%20in%20a%20linear%20regression\).](https://towardsdatascience.com/gradient-descent-algorithm-a-deep-dive-cf04e8115f21#:~:text=Gradient%20descent%20(GD)%20is%20an,e.g.%20in%20a%20linear%20regression).) [Accessed:26/02/2024].
- [10] ml4a. How neural networks are trained, 2020. https://ml4a.github.io/ml4a/how_neural_networks_are_trained/ [Accessed:27/02/2024].
- [11] mehran@mldawn.com. Back-propagation through cross-entropy softmax, 2021. <https://www.mldawn.com/back-propagation-with-cross-entropy-and-softmax/> [Accessed:29/02/2024].

