



Escuela Técnica Superior de
Ingeniería Informática

TRABAJO FIN DE GRADO

Detección Visual de Anomalías para Robots Móviles

Realizado por
Adrián Candil Salas

Para la obtención del título de
Grado en Ingeniería Informática - Tecnologías Informáticas

Dirigido por
Rocío González Díaz
Javier Perera Lago

En el departamento de
Matemática Aplicada I

Convocatoria de Julio, curso 2024

En primer lugar me gustaría no solo dedicar este trabajo, sino todo mi recorrido académico a mis padres Jaime y Luz María. Siempre han sido y serán un pilar fundamental en mi vida y sin ellos no habría llegado hasta aquí.

También me gustaría dedicarselo a mi hermmmana Lidia , para la cual me esfuerzo por ser un ejemplo a seguir.

Así mismo, quisiera también que esta dedicatoria sea dirigida hacia el resto de mi familia y a mis amigos que son para mí como hermanos. Especialmente a Nacho, Enrique y Juan Pablo los cuales han sido parte de este camino a diario durante mi carrera universitaria y han vivido tanto los buenos como los malos momentos conmigo y a Cristina, que a pesar de no ser nuestra compañera de piso ha vivido mucho con nosotros. Además deseo mencionar a mi compañero Jose con en cual he luchado muchas batallas en la facultad.

Finalmente agradecer a mi pareja Aida, la cual me ha dado el apoyo que necesitaba para no tropezar justo al final del camino.

Espero que ser parte de esta dedicatoria refleje el cariño y el agradecimiento que siento por todos los que me habéis acompañado. Sin vosotros no habría llegado aquí.

Agradecimientos

Quiero expresar mi más sincero agradecimiento a todas aquellas personas que contribuyeron de manera significativa al desarrollo de este trabajo.

En primer lugar, deseo agradecer a mis tutores Rocío González Díaz y Javier Perera Lago por su orientación y dedicación a lo largo del proyecto. Sus conocimientos y consejos fueron fundamentales para el desarrollo del mismo.

Una vez más quiero agradecer a mi familia y amigos, los cuales me han brindado siempre su cariño y apoyo. Para mí han sido y serán siempre cruciales.

Por último quiero reconocer a mis compañeros y profesores, quienes compartiendo conmigo sus conocimientos han contribuido a mi desarrollo académico.

A todos a los que habéis estado ahí siempre, de corazón os estoy eternamente agradecido.

Resumen

En este trabajo se realiza un estudio sobre una metodología aplicada al reconocimiento de anomalías mediante inteligencia artificial en las imágenes que toman las cámaras de robots móviles.

El documento presenta el desarrollo del estudio de esta metodología introduciendo previamente los conceptos necesarios para su entendimiento. He tratado de realizar el diseño y la implementación de varios de los métodos propuestos. El código se encuentra en GitHub. La arquitectura consta de dos partes combinadas que resultan en la clasificación de una imagen en normal o anómala.

Abstract

In this work, a study is carried out on a methodology applied to the recognition of anomalies using artificial intelligence in the images taken by mobile robot cameras.

The document presents the development of the study of this methodology, previously introducing the concepts necessary for its understanding. I have tried to carry out the design and implementation of several of the proposed methods. The code is located on GitHub. The architecture consists of two combined parts that result in the classification of an image as normal or anomalous.

Índice general

1	Introducción	1
2	Metodología	2
3	Planificación temporal y coste	3
4	Conceptos previos	6
4.1.	Red neuronal	6
4.1.1.	Perceptrón multicapa	7
4.2.	Convolución	8
4.3.	Red neuronal convolucional	10
4.4.	Autoencoder	12
4.5.	Entrenamiento de una red	12
4.5.1.	Dropout	13
4.5.2.	Aumento de datos (data augmentation)	13
4.5.3.	Regularización L1 y L2	14
4.5.4.	Early stopping, ReduceLROnPlateau y gradient clipping	14
4.5.5.	Reducción de las dimensiones de la red	15
4.5.6.	Validación cruzada	15
5	Ejecución del trabajo	16
5.1.	Conjunto de datos	16
5.2.	Estudio y definición de la arquitectura objetivo	17
5.2.1.	Autoencoder como detector de anomalías	18
5.2.2.	Autoencoder + MLP	18
5.3.	Decisión del entorno de desarrollo	19
5.4.	Desarrollo del código	20
5.4.1.	Importación de los datos a GoogleColab	20
5.4.2.	Versiones y entrenamiento de los modelos	21
5.4.3.	Ejemplo de uso del autoencoder	28
5.4.4.	Autoencoder como detector de anomalías	29
5.4.5.	Codificación de las imágenes	32
5.4.6.	Perceptrón multicapa	32
6	Conclusiones	33
	Bibliografía	34

Índice de figuras

4.1. Red neuronal	6
4.2. Neurona	7
4.3. Perceptrón multicapa	7
4.4. Operación de convolución	8
4.5. Convolución paso a paso	8
4.6. Convolución 3D	9
4.7. Feature map 3D	9
4.8. Capa convolucional	10
4.9. Red convolucional + red neuronal	10
4.10. Pooling vs Upsample	11
4.11. Autoencoder	12
5.1. Ejemplos de imágenes	16
5.2. Estructura autoencoder paper, [1, Página 5]	18
5.3. Reconstrucción sin anomalía	28
5.4. Reconstrucción con anomalía	29
5.5. Ejemplos error cuadrático medio	30

Índice de tablas

3.1. Tiempos estimados y reales	3
3.2. Tiempo estimado vs tiempo real	4
3.3. Coste de personal	4
3.4. Coste de servicios	5
3.5. Coste total	5
5.1. Dataset	17

Índice de extractos de código

5.1. Perceptrón multicapa	19
5.2. Código para importar las imágenes	20
5.3. Resize de las imágenes	21
5.4. Autoencoderv1	22
5.5. Autoencoder v2	23
5.6. Autoencoder v3	24
5.7. Data augmentation	25
5.8. Cross Validation	25
5.9. Autoencoder v4	26
5.10. Código para importar las imágenes	27
5.11. Reconstrucción sin anomalía	28
5.12. Reconstrucción con anomalía	28
5.13. Error cuadrático medio	29
5.14. Clasificador autoencoder	30
5.15. Clasificador autoencoder por batches	31
5.16. Codificación de las imágenes	32

1. Introducción

Día tras día, las inteligencias artificiales (IAs) están tomando un rol cada vez mas importante dentro de la sociedad, extendiendose en todos los sectores. Desde chatboxes como chat-GPT o IAs que generan imágenes hasta el desarrollo de coches o robots autónomos.

En este trabajo vamos a centrarnos en cubrir una necesidad que surge en el desarrollo de estos últimos: ¿Cómo puede un robot equipado con una cámara detectar anomalías?

Para responder a esta respuesta centraremos nuestra atención en la arquitectura propuesta por el grupo de investigadores formado por Dario Mantegazza , Alessandro Giusti , Luca Maria Gambardella, and Jérôme Guzzi. Esta consiste en primero reducir la dimensión de la imagen tomada por el robot y posteriormente aplicar diferentes modelos de clasificación sobre esta reducción para poder decidir si la imagen contiene o no una anomalía.

Para el estudio de esta metodología, además de basarnos en el artículo que la desarrolla, se han revisado numerosas fuentes bibliográficas y sobre todo se ha tratado de simular los experimentos llevados a cabo por el equipo mencionado previamente. Este código queda disponible en mi repositorio de [github](#).

Este trabajo se divide en varias partes. La sección [2](#), que contempla la metodología seguida durante el desarrollo del trabajo. A continuación encontramos la sección [3](#), la cual describe la planificación inicial del proyecto así como el desarrollo real de este. Después en la sección [4](#) hablamos de los conceptos previos necesarios para el entendimiento de la arquitectura desarrollada. La sección [5](#) explica la arquitectura propuesta por el paper de estudio [\[1\]](#) y explica la implementación realizada por el alumno. Finalmente en la sección [6](#) comentamos brevemente las conclusiones del trabajo.

2. Metodología

La realización de este trabajo ha conestado de un trabajo paralelo de desarrollo de código y de investigación.

Comencé estudiando el artículo *An Outlier Exposure Approach to Improve Visual Anomaly Detection Performance for Mobile Robots* [1] de los investigadores mencionados previamente en la introducción. Una vez tuve un nivel de comprensión suficiente, comencé a desarrollar el código. El desarrollo del código ha sido la parte a la que mas tiempo he dedicado como veremos mas adelante en la sección 3. Este desarrollo se basa en la traducción del código aportado en el paper mencionado previamente, que está en PyTorch [2], a Tensorflow [3] en un notebook de GoogleColab [4]. Pytorch y Tensorflow son las librerías mas desarrolladas para la implementación de modelos de inteligencia artificial. Por otra parte GoogleColab es una herramienta que ofrece Google para trabajar con notebooks de python. Debido a las limitaciones de GoogleColab, hemos tenido que encontrar una gran cantidad de soluciones, que veremos durante el documento, para poder llevar a cabo los experimentos.

Como gran parte de los experimentos han sido largos entrenamientos de modelos de IA, durante la supervisión de los entrenamientos realicé gran parte de la investigación del trabajo.

A lo largo del desarrollo del trabajo he realizado reuniones cada dos semanas con Rocío González Díaz y Javier Perera Lago, tutores de este TFG.

3. Planificación temporal y coste

En esta sección vamos a ver la planificación del proyecto así como el coste estimado si fuese un proyecto real. La planificación expone las actividades realizadas, los tiempos estimados de estos y el tiempo real que hemos dedicado a cada actividad. La planificación ha ido variando durante el proyecto como veremos a continuación.

Por otra parte encontramos el análisis del coste que supondría este proyecto si fuese un proyecto financiado. Además también detallará el coste real de este, ya que se han tenido que adquirir ciertos servicios.

3.1 Planificación

La idea inicial era entregar en la primera convocatoria. Por desgracia debido a la imposibilidad de hacerlo, finalmente la entrega se realizó en la segunda un mes después. Un proyecto se divide en cuatro fases:

1. Iniciación
2. Planificación
3. Ejecución
4. Seguimiento y Control

La fase de iniciación es la fase inicial de un proyecto donde se identifican las necesidades y se evalúa la viabilidad del este. Durante la fase de planificación se establece un plan concreto para llevar a cabo el desarrollo del proyecto. La fase de ejecución consiste en llevar a cabo el plan establecido en la planificación. Finalmente, el seguimiento y control consiste en la monitorización de la ejecución, asegurando que se está llevando a cabo de forma correcta.

La siguiente tabla muestra el desarrollo de las siguientes fases.

	Iniciación	Planificación	Ejecución	Seguimiento y Control
Inicio	13/09/2023	03/10/2023	06/11/2023	03/10/2023
Fin estimado	03/10/2023	10/10/2023	23/05/2024	17/05/2024
Fin real	03/10/2023	17/10/2023	18/06/2024	14/06/2024

Tabla 3.1: Tiempos estimados y reales

Las actividades realizadas a lo largo del trabajo han sido las siguientes: planificación del trabajo, reuniones de seguimiento, desarrollo de código, estudio de la bibliografía, presentación, desarrollo de la memoria.

A lo largo del desarrollo del proyecto he tenido que dedicar mucho del tiempo de estudio de la bibliografía y dedicarlo al desarrollo de código aunque como ya mencioné en la metodología parte de este tiempo de desarrollo fue supervisar que los entrenamientos fuesen como se esperaba y mientras tanto investigaba. Las horas reales quedarían así:

Tarea	Horas estimadas	Horas reales
Planificación	6:00h	6:00h
Reuniones	7:00h	8:00h
Desarrollo del código	150:00h	203:00h
Estudio de la bibliografía	100:00h	55:00h
Desarrollo de la presentación	7:00h	6:30h
Desarrollo de la memoria	30:00h	40:00h
Total	300:00h	318:30h

Tabla 3.2: Tiempo estimado vs tiempo real

3.2 Coste

Para el analisis del coste que supondría el proyecto consideraremos tanto el pago al personal, como los gastos de servicios externos.

Comenzamos pues con el coste de personal. Vamos a asumir que el alumno ya es considerado un Ingeniero informático y por tanto su sueldo sería el correspondiente [5].

Rol	Sueldo	Horas	Coste
Ingeniero informático	13.85€/h	318:30h	4411.25€

Tabla 3.3: Coste de personal

Vamos a considerar ahora los servicios utilizados [3.4](#). Consideramos en primer lugar los servicios básicos que son la luz consumida y el acceso a internet. Por otra parte se ha tenido que contratar varios servicios de Google: Google One, Unidades de ejecución en GoogleColab y ColabPro. Los servicios de GoogleOne nos proporcionan un mayor almacenamiento en GoogleDrive, lo cual nos dará la capacidad de almacenar los datasets ahí. Por su parte ColabPro nos dará máquinas más potentes que habilitarán la ejecución del código. Para poder usar estas máquinas Colab consume unidades de ejecución. Dependiendo de cuantos recursos consumes se restará un número aproximado de unidades de ejecución. Aproximadamente unas 4 por hora. Cada adquisición nos proporciona una 100 unidades y otras 100 al mes con la subscipción a ColabPro. El gasto de internet corresponde al pago mensual aproximado que realiza el alumno.

Servicio	Precio	Meses	Total
Luz	20€/mes	8	160€
Internet	11€/mes	8	88€
Google One	1.99€/mes	3	5.97€
Unidades de ejecución	11.99€/adquisición	3	35.97€
Colab Pro	11.99€/mes	2	23.98€
Total			313.92€

Tabla 3.4: Coste de servicios

El coste total del proyecto será por tanto la suma del coste de personal y de servicios.

-	Coste
Coste de personal	4411.25€
Coste de servicios	313.92€
Total	4713.18€

Tabla 3.5: Coste total

4. Conceptos previos

Para poder comprender en su totalidad la arquitectura trabajada a lo largo de este proyecto, primero debemos explicar una serie de conceptos. Concretamente, realizaremos una breve introducción a las redes neuronales, definiremos el perceptrón multicapa, describiremos la operación de convolución, estableceremos qué es una red neuronal convolucional y especificaremos concretamente cómo está estructurado un autoencoder y finalmente hablaremos del entrenamiento de modelos y de las diferentes técnicas utilizadas para evitar el sobreajuste.

4.1. Red neuronal

Una red neuronal es una estructura matemática inspirada en el sistema nervioso. Funciona utilizando un elevado número de unidades de procesamiento que tratan de simular a las neuronas. Estas neuronas se interconectan formando la estructura de la red. [6, Página 58]

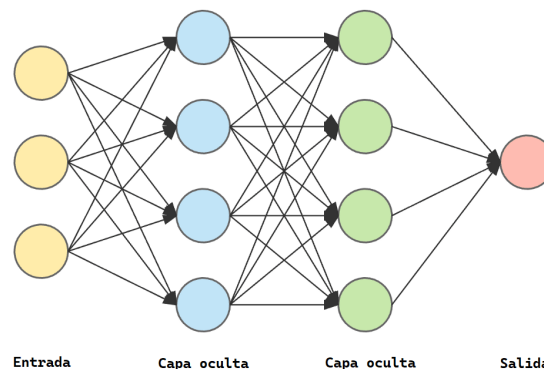


Figura 4.1: Red neuronal

La estructura de una red neuronal como podemos observar en la Figura 4.1 consta de una capa de entrada, la cual recibe directamente la información que vamos a procesar; capas ocultas, las cuales contienen las neuronas cuya distribución e interconexión determina la topología de la red y finalmente, una capa de salida que, dependiendo de la función de nuestra red proporcionará un tipo de información u otra (p.e: una clasificación, una imagen modificada, etc.)

En cada neurona se realiza la siguiente operación llamada sinapsis:

$$a_j = g \left(\sum_{i=0}^n w_{ij} a_i \right)$$

Donde a_j es la activación de la neurona j , g es la función de activación común en toda la capa, w_{ij} es el peso asociado a la conexión entre la neurona i de la capa anterior y la neurona j , y a_i son las activaciones provenientes de las neuronas en la capa anterior. La función de activación es una función que se aplica al resultado del computo de la sinapsis para introducir la no linealidad en el modelo, permitiendo que este aprenda y represente relaciones complejas de datos.

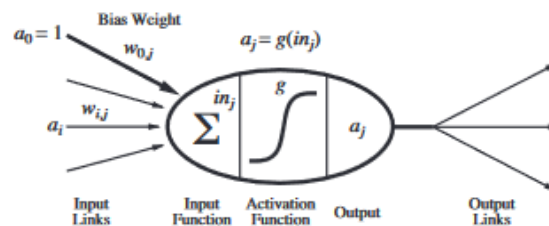


Figura 4.2: Neurona

4.1.1. Perceptrón multicapa

Un perceptrón multicapa es la forma más simple de red neuronal y la base de los modelos más avanzados desarrollados en el aprendizaje profundo. Normalmente, se utilizan en problemas de clasificación donde es necesario asignar etiquetas a observaciones (binarias o multinomiales) basadas en los datos de entrada.

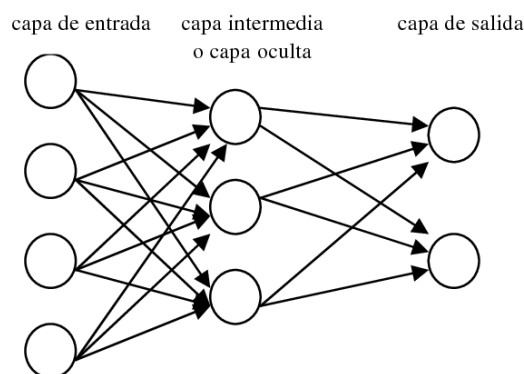


Figura 4.3: Perceptrón multicapa

4.2. Convolución

Una convolución es una operación matemática que combina dos funciones para producir una tercera función. Esta operación se utiliza para determinar cómo la forma de una función se modifica por otra. Sin embargo, la operación que nos interesa es la convolución desde un punto de vista matricial. [6, Página 122]

En el tratamiento de imágenes, estas son tratadas como matrices. Vamos a centrarnos por tanto en la convolución matricial, operación que, combinada con una red neuronal da lugar a una red neuronal convolucional. [6, Página 36]

La operación de convolución matricial se basa en hacer un barrido de la imagen con un kernel o filtro dando como resultado otra matriz del mismo tamaño o de tamaño reducido llamada mapa de características (feature map). [6, Página 122-127]

Veamos cómo se aplica la operación con un ejemplo. Tenemos una imagen de entrada 5x5 y un filtro 3x3. En este caso concreto el mapa de características tendrá un tamaño reducido tomando como referencia la imagen de entrada.

<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>2</td><td>9</td><td>11</td><td>22</td><td>42</td></tr> <tr><td>12</td><td>43</td><td>2</td><td>0</td><td>65</td></tr> <tr><td>9</td><td>87</td><td>65</td><td>2</td><td>90</td></tr> <tr><td>21</td><td>21</td><td>44</td><td>7</td><td>2</td></tr> <tr><td>45</td><td>2</td><td>21</td><td>2</td><td>0</td></tr> </table>	2	9	11	22	42	12	43	2	0	65	9	87	65	2	90	21	21	44	7	2	45	2	21	2	0	⊗	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>0</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>0</td></tr> <tr><td>0</td><td>0</td><td>1</td></tr> </table>	0	0	1	1	0	0	0	0	1	=	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>88</td><td>67</td><td>134</td></tr> <tr><td>55</td><td>94</td><td>132</td></tr> <tr><td>107</td><td>25</td><td>134</td></tr> </table>	88	67	134	55	94	132	107	25	134
2	9	11	22	42																																											
12	43	2	0	65																																											
9	87	65	2	90																																											
21	21	44	7	2																																											
45	2	21	2	0																																											
0	0	1																																													
1	0	0																																													
0	0	1																																													
88	67	134																																													
55	94	132																																													
107	25	134																																													
Matriz imagen Entrada a capa convolucional 5x5		Kernel o filtro 3x3		Feature map Salida de la capa convolucional 3x3																																											

Figura 4.4: Operación de convolución

El feature map es obtenido paso a paso de la siguiente forma definida en 4.5. El kernel se va desplazando por la imagen y va realizando una multiplicación element-wise con la submatriz que superpone y suma los resultados para ir obteniendo cada uno de los valores del feature map.

<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>2₀</td><td>9₀</td><td>11₁</td><td>22</td><td>42</td></tr> <tr><td>12₁</td><td>43₀</td><td>2₀</td><td>0</td><td>65</td></tr> <tr><td>9₀</td><td>87₀</td><td>65₁</td><td>2</td><td>90</td></tr> <tr><td>21</td><td>21</td><td>44</td><td>7</td><td>2</td></tr> <tr><td>45</td><td>2</td><td>21</td><td>2</td><td>0</td></tr> </table>	2 ₀	9 ₀	11 ₁	22	42	12 ₁	43 ₀	2 ₀	0	65	9 ₀	87 ₀	65 ₁	2	90	21	21	44	7	2	45	2	21	2	0	=	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>88</td><td></td><td></td></tr> <tr><td></td><td></td><td></td></tr> <tr><td></td><td></td><td></td></tr> </table>	88								
2 ₀	9 ₀	11 ₁	22	42																																
12 ₁	43 ₀	2 ₀	0	65																																
9 ₀	87 ₀	65 ₁	2	90																																
21	21	44	7	2																																
45	2	21	2	0																																
88																																				
convolución		convolución																																		
<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>2</td><td>9</td><td>11</td><td>22</td><td>42</td></tr> <tr><td>12</td><td>43</td><td>2</td><td>0</td><td>65</td></tr> <tr><td>9</td><td>87</td><td>65</td><td>2</td><td>90</td></tr> <tr><td>21</td><td>21</td><td>44</td><td>7</td><td>2</td></tr> <tr><td>45</td><td>2</td><td>21</td><td>2</td><td>0</td></tr> </table>	2	9	11	22	42	12	43	2	0	65	9	87	65	2	90	21	21	44	7	2	45	2	21	2	0	=	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>88</td><td></td><td></td></tr> <tr><td>55</td><td></td><td></td></tr> <tr><td>107</td><td></td><td></td></tr> </table>	88			55			107		
2	9	11	22	42																																
12	43	2	0	65																																
9	87	65	2	90																																
21	21	44	7	2																																
45	2	21	2	0																																
88																																				
55																																				
107																																				
convolución		convolución																																		
<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>2</td><td>9</td><td>11</td><td>22</td><td>42</td></tr> <tr><td>12</td><td>43</td><td>2</td><td>0</td><td>65</td></tr> <tr><td>9</td><td>87</td><td>65</td><td>2</td><td>90</td></tr> <tr><td>21</td><td>21</td><td>44</td><td>7</td><td>2</td></tr> <tr><td>45</td><td>2</td><td>21</td><td>2</td><td>0</td></tr> </table>	2	9	11	22	42	12	43	2	0	65	9	87	65	2	90	21	21	44	7	2	45	2	21	2	0	=	<table border="1" style="border-collapse: collapse; text-align: center;"> <tr><td>88</td><td>67</td><td></td></tr> <tr><td>55</td><td></td><td></td></tr> <tr><td>107</td><td></td><td></td></tr> </table>	88	67		55			107		
2	9	11	22	42																																
12	43	2	0	65																																
9	87	65	2	90																																
21	21	44	7	2																																
45	2	21	2	0																																
88	67																																			
55																																				
107																																				
convolución		convolución																																		
...		...																																		

Figura 4.5: Convolución paso a paso

En la primera multiplicación de 4.5, la operación sería la siguiente:

$$2 \cdot 0 + 9 \cdot 0 + 11 \cdot 1 + 12 \cdot 1 + 43 \cdot 0 + 2 \cdot 0 + 9 \cdot 0 + 87 \cdot 0 + 65 \cdot 1 = 88$$

El caso que acabamos de ver es válido para imágenes en blanco y negro que están formadas por un solo canal. Llamamos canal en este contexto a cada una de las dimensiones que contiene información relevante de un dato de entrada. Sin embargo esta operación puede realizarse también en imágenes a color que están formadas por 3 canales. Esto es debido a que las imágenes en color están formadas por los canales RGB. Al representarla de forma matricial, una imagen es el tensor de cada una de las matrices correspondientes a cada uno de los canales del RGB. De esta forma una imagen será un tensor 3D y por tanto el kernel deberá tener también tres canales.

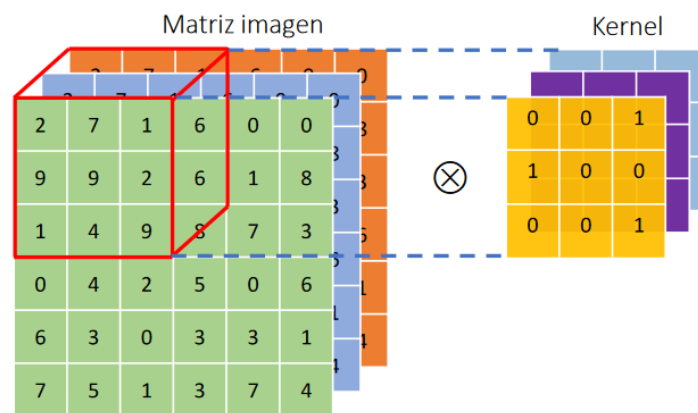


Figura 4.6: Convolución 3D

El feature map se obtendría en este caso de la siguiente manera definida en 4.7:

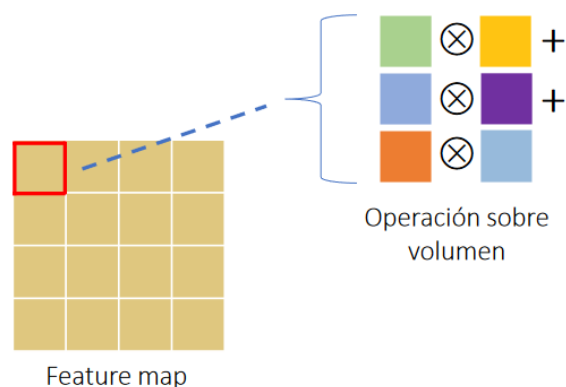


Figura 4.7: Feature map 3D

De esta forma por cada kernel extrae un feature map unidimensional realizando la suma de la convolución de cada canal.

4.3. Red neuronal convolucional

Mezclando las redes neuronales y la convolución obtenemos uno de los instrumentos mas populares de los últimos tiempos: la red neuronal convolucional. [6, Página 30]

Estas redes son redes neuronales que en vez de aprender los pesos que interconectan a las neuronas, la red aprenderá los pesos de los k filtros aplicados a la entrada de la capa, generando los k feature maps correspondientes y aplicando la función de activación correspondiente a la capa para dar el output. Todo este proceso conforma una capa convolucional.[6, Página 120-177]

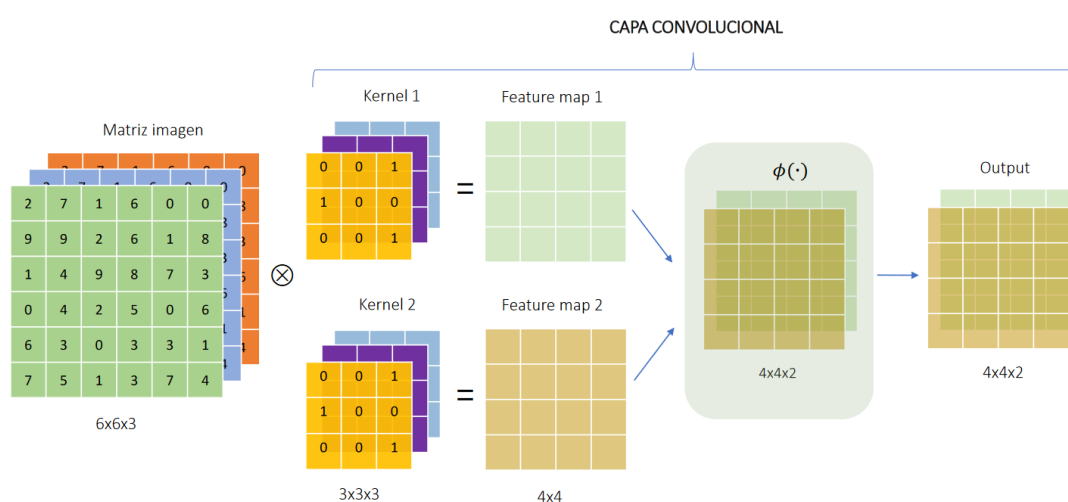


Figura 4.8: Capa convolucional

La sucesión de diferentes capas convolucionales, distribuidas de diferentes maneras dependiendo de la funcionalidad que queramos que tenga la red resultará en una red neuronal convolucional.

Normalmente se concatenan varias capas convolucionales para reducir la dimensión y aprender las características de una imagen y se concatena a un perceptrón multicapa que es el que realiza la clasificación. Un ejemplo de esta es el siguiente 4.9, el cual encontramos en [6, Página 153]:

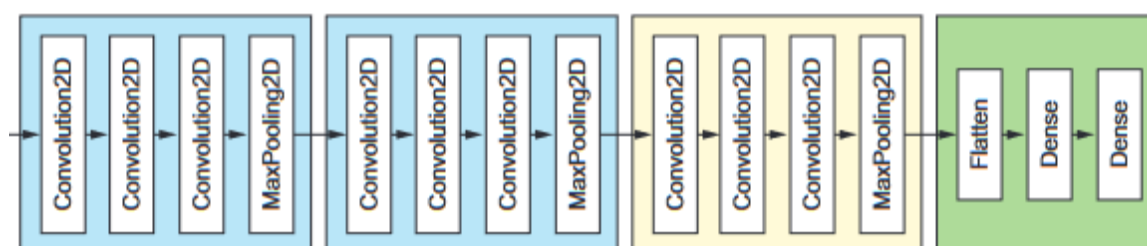


Figura 4.9: Red convolucional + red neuronal

Concretamente para construir nuestro modelo como veremos más adelante hemos utilizado capas Conv2D y Conv2DTranspose y en la versión inicial de nuestro modelo 5.4, MaxPooling2D y Upsample2D.

Conv2D realiza la operación de convolución descrita hasta el momento. Por otra parte las capas Conv2DTranspose, realiza la operación llamada convolución transpuesta o deconvolución. Esta operación consiste en aumentar la dimensión espacial y reconstituir la imagen. Teniendo en cuenta que Conv2D reduce la dimensión y extrae características y Conv2DTranspose aumenta la dimensión y reconstituye la imagen, podemos considerarlas operaciones opuestas.

Vamos a definir brevemente la operación de MaxPooling2D y Upsample2D que utilizamos en versiones iniciales del código.

La operación de MaxPooling consiste en resumir las características aplicando un kernel a los feature maps resultantes de una convolución aplicando la función máximo.

Por otro lado la operación de upsample es la opuesta a la de pooling, aumentando en el proceso la dimensión espacial de un mapa de características más pequeño mediante interpolación.

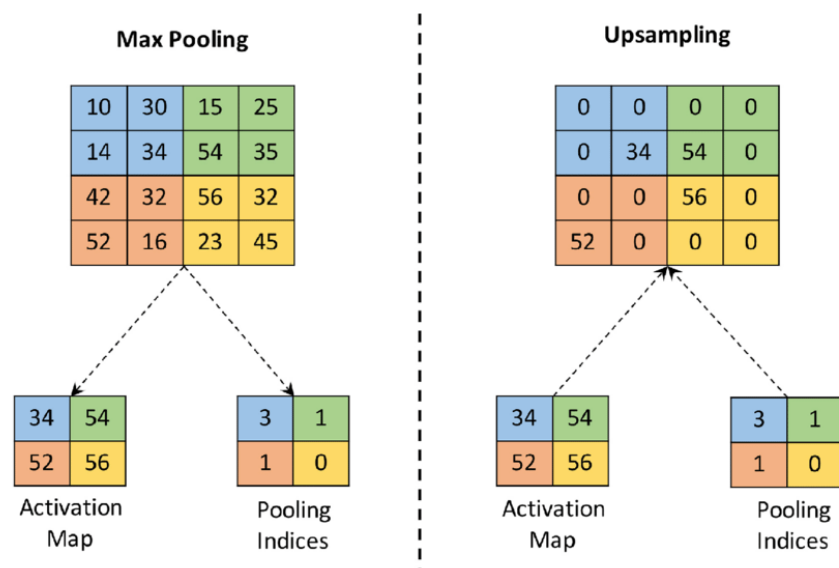


Figura 4.10: Pooling vs Upsample

En nuestro estudio, vamos a profundizar concretamente en la estructura del autoencoder.

4.4. Autoencoder

Un autoencoder es una red neuronal convolucional formada por 3 bloques: el codificador, el cuello de botella y el decodificador. El objetivo del autoencoder es reducir la dimension de la entrada y luego reconstruirla minimizando el error cometido durante el proceso.[7, Página 1] Esta técnica se usa por ejemplo con el fin de reducir el peso de almacenamiento de los datos de entrada o como en nuestro caso, para preprocesar estos datos y dárselos como entrada a otros modelos de entrenamiento.

El codificador recibe un dato de entrada, para nuestro caso será una imagen, y reducirá su dimensión. Nuestro codificador recibirá imagenes 64x64x3 y las reduce a embeddings unidimensionales de tamaño 128.

El decodificador se encarga de reconstruir la codificación producida por el encoder. En nuestro caso reconstruirá el embedding 128 a una imagen 64x64x3 de nuevo.

El cuello de botella es la unión entre el codificador y el decodificador y donde encontramos los elementos codificados.

Para visualizar mejor esta red neuronal podemos imaginarnos una pajarita ya que su forma es identica a estas. En 4.11 podemos apreciar la forma mencionada.

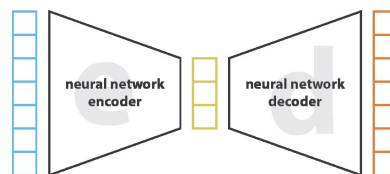


Figura 4.11: Autoencoder

4.5. Entrenamiento de una red

Como hemos mencionado previamente, las diferentes redes tienen que ser entrenadas sobre un conjunto de datos concreto para que se ajusten a una tarea. En nuestro caso, trabajamos con imágenes así que la arquitectura entrenará sobre un conjunto de imágenes.

La idea es que las imágenes se introduzcan en el autoencoder y este aprenda a reconstruirlas reduciendo el valor de la función de pérdida. Esta función trata de cuantificar la diferencia de los resultados que da la red y los valores reales de los datos de entrenamiento. Dependiendo del propósito de la red se usarán diferentes funciones para esta finalidad. En nuestro caso usaremos el error medio cuadrático el entrenamiento de ul autoencoder y la entropía binaria cuadrada para el entrenamiento de un perceptrón multicapa.

$$\text{Error cuadrático medio} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{Entropía binaria cruzada} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

El aprendizaje consistirá pues en que los pesos de los filtros de las diferentes capas del autoencoder se ajusten al conjunto de entrenamiento que le hemos proporcionado. Sin embargo debido al alto número de parámetros que se utilizan, enfrentamos el problema del sobreajuste. En cada epoch del entrenamiento los pesos se van actualizando a medida que el valor de la función de pérdida baja. La actualización de los pesos se hace mediante el algoritmo de descenso por el gradiente en el que en cada epoch cada peso se actualiza de la siguiente manera:

$$w_{ij} \leftarrow w_{ij} - \eta \frac{\partial L}{\partial w_{ij}}$$

Donde w_{ij} es el valor del peso, η es el factor de aprendizaje, un valor preestablecido que actúa de regulador de la actualización de los pesos y finalmente $\frac{\partial L}{\partial w_{ij}}$ es el gradiente de la función L respecto a w_{ij}

El sobreajuste es a efectos prácticos que la red se ajuste demasiado a los datos de entrada y no generalice bien sobre otros datos. Para evitar esto existen múltiples técnicas que hemos aplicado durante los experimentos.

4.5.1. Dropout

El dropout consiste en durante el entrenamiento apagar el entrenamiento de un porcentaje de los pesos aleatoriamente para que así no aprendan características específicas de conjunto de entrenamiento. En cada epoch se apagan diferentes pesos. [8, Página 2]

4.5.2. Aumento de datos (data augmentation)

La técnica del aumento de datos se basa en realizar diferentes modificaciones a las imágenes de entrada para que manteniendo una base similar, tengamos una cantidad de datos más amplia con diferentes características. De esta forma la red es capaz de aprender las características deseadas siendo capaz de generalizar mejor sobre otros datos. [9, Página 7]

Las modificaciones más comunes son flips horizontales o verticales, rotaciones de diferentes grados, zoom-in o zoom-out, modificación en el brillo de la imagen o la introducción de ruido de sal y pimienta. [9, Páginas 7-10]

4.5.3. Regularización L1 y L2

La regularización de pesos se basa en añadir una penalización en la función de pérdida. [8, Páginas 2-3]

- La función de pérdida con regularización L1 se define como:

$$\text{Loss} = \text{Loss}_{\text{original}} + \lambda \sum_i |w_i|$$

La regularización L1 agrega una penalización proporcional a la suma de los valores absolutos de los coeficientes (pesos) del modelo. Esta técnica se conoce como Lasso (Least Absolute Shrinkage and Selection Operator).

-La función de pérdida con regularización L2 se define como:

$$\text{Loss} = \text{Loss}_{\text{original}} + \lambda \sum_i w_i^2$$

La regularización L2 agrega una penalización proporcional a la suma de los cuadrados de los coeficientes (pesos) del modelo. Esta técnica se conoce como Ridge Regression.

De esta manera los pesos van cada vez aprendiendo menos para así impedir que se sobreaprendan características específicas del conjunto de entrenamiento.

4.5.4. Early stopping, ReduceLRonPlateau y gradient clipping

Otra práctica común para evitar el sobreajuste es añadir callbacks al entrenamiento. Un callback basicamente es una función que puede aplicarse durante el entrenamiento para personalizar y controlar su comportamiento.

En este caso usaremos early stopping que como su propio nombre indica consiste en parar el entrenamiento antes de que se complete al cien por cien. Para ello establecemos que parámetro estará supervisando el callback y si se da un empeoramiento durante un número de epochs determinado (al que llamamos paciencia) el entrenamiento para, normalmente restaurando los pesos que ofrecen un mejor resultado.

Por otra parte ReduceLRonPlateau es un callback que al igual que early stopping irá monitoreando un parámetro preestablecido. Sin embargo, en vez de parar el entrenamiento, cuando el parámetro empeora n epochs seguidos, lo que hará ReduceLRonplateau, será reducir el learning rate en un factor preestablecido.

Por último el gradient clipping es una técnica que consiste en limitar el valor máximo de los durante el entrenamiento.

4.5.5. Reducción de las dimensiones de la red

El último método que hemos utilizado para tratar de combatir el sobreajuste es reducir las dimensiones de la red. Al eliminar capas, tenemos menos parámetros lo cual beneficia de manera directa la reducción del sobreajuste ya que reduce la capacidad de aprender patrones complejos. [10, Página 2]

4.5.6. Validación cruzada

La validación cruzada es una técnica que consiste en dividir el conjunto de entrenamiento en k subconjuntos. Una vez realizada la división, se utilizan $k-1$ conjuntos para entrenamiento y el restante como conjunto de validación. La distribución se va permutando k iteraciones hasta haber realizado todos los subentrenamientos. De esta manera se consigue reducir el sobreajuste.

Una vez explicados estos conceptos podemos proceder a desarrollar la explicación del trabajo realizado.

5. Ejecución del trabajo

Una vez introducidos los conceptos previos, vamos a pasar a desarrollar la explicación del trabajo realizado.

Empecemos definiendo cuál es el objetivo de este trabajo. La idea es, utilizando como referencia la arquitectura presentada en el artículo *An Outlier Exposure Approach to Improve Visual Anomaly Detection Performance for Mobile Robots* [1], hacer una traducción de esta, que originalmente está implementada en PyTorch, a Tensorflow y ver si podían mejorar los resultados de esta manera.

A continuación vamos a explicar todo el proceso que ha conllevado la traducción del código. Desde el estudio de la arquitectura, pasando por el conjunto de decisiones que se han tenido que ir tomando sobre la marcha para solventar las diferentes complicaciones que iban surgiendo y finalmente expondremos los resultados obtenidos.

5.1. Conjunto de datos

Para poder entender tanto la arquitectura como los problemas encontrados durante el desarrollo tenemos que explorar el conjunto de datos.

Este esta formado por 132838 imágenes 512x512 RGB, 105122 normales y 27716 con anomalías [1, Página 4, apartado A]. El dataset se puede encontrar en <https://zenodo.org/records/7074958>.

Ejemplos de estas imágenes son las siguientes 5.1:



Figura 5.1: Ejemplos de imágenes

Las imágenes están divididas de la siguiente manera [5.1](#):

-	Training set	Validation set	Testing set
Normal frames	56066	2071	46585
Anomalous frames	18201	-	9515

Tabla 5.1: Dataset

Entrenaremos el autoencoder con el conjunto de imágenes normales del conjunto de entrenamiento.

Una vez que el autoencoder esté entrenado lo usaremos de las siguientes formas:

1. Como detector de anomalías.
2. Reductor de dimensiones para generar un conjunto de entrenamiento para un perceptrón multicapa.

Vamos ahora a describir con un poco más de profundidad ambas arquitecturas objetivo.

5.2. Estudio y definición de la arquitectura objetivo

Una vez tuvimos una comprensión inicial del paper [\[1\]](#), definimos el que sería el objetivo del trabajo: realizar la implementación de un autoencoder capaz de reducir la dimensionalidad de las imágenes tomadas por un robot y usar esta reducción como conjunto de entrenamiento para modelos más simples como es el perceptrón multicapa. La utilidad de esto es conseguir una estructura capaz de detectar si en la imagen hay anomalías prácticamente en tiempo real. Además el autoencoder puede ser usado de por sí solo como clasificador de anomalías como veremos más adelante.

Vamos a hacer ahora una definición más extensa de los dos modelos hemos implementado, propuestos en [\[1\]](#). Por un lado hemos desarrollado el autoencoder como detector de anomalías y por otro hemos combinado el autoencoder, usandolo como reductor de la dimensión los elementos del conjunto de entrenamiento, con un perceptrón multicapa que funcionará como clasificador dando así pie a dos arquitecturas distintas que cumplen con la misma función.

1. Autoencoder
2. Autoencoder + MLP

5.2.1. Autoencoder como detector de anomalías

Como ya hemos mencionado previamente en este trabajo, un autoencoder codifica y luego reconstruye imágenes reduciendo el error cuadrático medio [1, Página 3]. Podemos usar este error en la reconstrucción como puntuación de anomalía de la siguiente manera:

1. Entrenamos el autoencoder con un conjunto de datos en el que las imágenes estén libres de anomalías.
2. Cuando pasamos una imagen por la red si la reconstruye con un valor bajo en la función de pérdida tendremos una imagen sin anomalías. Por otra parte si el valor es alto significa que la imagen original dista mucho de la reconstruida ya que el autoencoder no es capaz de reconstruir anomalías .

La estructura propuesta por los investigadores es la siguiente:



Figura 5.2: Estructura autoencoder paper, [1, Página 5]

Se usa LeakyRelu 5.1 como función de activación en todas sus capas menos en las que específicamente usa una función lineal. Después de cada convolución el tamaño de la imagen queda reducido a la mitad y en cada upsample queda multiplicado por dos [1, Página 5].

$$f(x) = \begin{cases} x & \text{si } x \geq 0 \\ \alpha x & \text{si } x < 0 \end{cases} \quad (5.1)$$

5.2.2. Autoencoder + MLP

Por otra parte utilizaremos el autoencoder entrenado para reducir la dimensionalidad del conjunto de entrenamiento que utilizamos para entrenar un perceptrón multicapa muy simple que funcionará como clasificador.

De esta manera las imágenes 64x64x3 quedan reducidas a un vector 128 haciendo que su procesamiento sea mucho menos costoso.

El conjunto de entrenamiento del perceptrón está formado por imágenes codificadas que contienen anomalías y otras que no lo hacen [1, Página 5, apartado D].

El perceptrón propuesto consta de tres capas de 256, 64 y 1 neuronas respectivamente con activaciones relu en las dos primeras y sigmoide en la última capa y usa como función de pérdida la entropía binaria cruzada.

Al ser un modelo tan simple se ha implementado sin ninguna modificación en nuestro código.

```
1 MLP = Sequential([
2     Dense(256, input_shape=(128,), activation='relu'),
3     Dense(64, activation='relu'),
4     Dense(1, activation='sigmoid')
5 ])
6
7 opt = keras.optimizers.Adam(learning_rate=0.0001)
8 MLP.compile(optimizer=opt, loss='binary_crossentropy', metrics=['accuracy'])
9 MLP.summary()
```

Extracto de código 5.1: Perceptrón multicapa

5.3. Decisión del entorno de desarrollo

Para desarrollar el código primero tuvimos que decidir en que entorno íbamos a desarrollarlo. En una de las primeras reuniones se decidió que se iba a generar un notebook de Python en el cual se realizaría el trabajo. Por tanto, esto nos dejó con dos opciones principalmente:

1. Trabajar en local con JupyterNotebook.
2. Usar GoogleColab, que como se ve en el apartado de planificación 3, fue la opción tomada.

La decisión de usar colab y no trabajar en local fue tomada ya que aportaba la ventaja de que no iba a ser el equipo del alumno el que trabajase los entrenamientos, quitándole una gran carga a este. Además nos solucionaba el trabajo de instalar las diferentes librerías utilizadas, ya que se pueden usar directamente en la herramienta de Google.

5.4. Desarrollo del código

En este apartado vamos a explicar cómo hemos llevado a cabo el desarrollo del código de la traducción a Tensorflow. El código está disponible en el siguiente [github](#). Estará dividido en los siguientes subapartados:

1. Importación de los datos a GoogleColab.
2. Versiones y entrenamiento de los autoencoders.
3. Uso del autoencoder como detector de anomalías.
4. Codificación de imágenes.
5. Optimización de la importación de los datos.
6. Perceptrón multicapa.

Durante la explicación de ambos subapartados iremos viendo las diferentes dificultades que hemos enfrentado.

5.4.1. Importación de los datos a GoogleColab

El primer problema que enfrentamos a la hora de llevar los datos a la nube fue que nos equivocamos de dataset. En el [github](#) de los investigadores hay diferentes dataset correspondientes a diferentes papers. Debido a una mala lectura de este, inicialmente usamos un dataset de otro paper que a pesar de usar imágenes similares a las de nuestro verdadero dataset, era mucho mas grande pesando 20GB aproximadamente. Esto hacía que subirlo al entorno de ejecución cada vez que fuésemos a trabajar no fuese una opción, por tanto solo quedaba utilizar el enlace entre Google Drive y GoogleColab. Adquirimos pues GoogleOne para poder tener un mayor almacenamiento disponible en la nube y procedimos a subir los datos. Comentar que debían estar comprimidos en zip para que la importación a Colab fuese mas eficiente.

```
1 extract_path = '/content/'
2
3 data = '/content/drive/MyDrive/Corridors v2.0.zip'
4
5 weighs = '/content/drive/MyDrive/autoencoder_weights_my_version_1st_train.zip'
6
7 with zipfile.ZipFile(data, 'r') as zip_ref:
8     zip_ref.extractall(extract_path)
9
10 with zipfile.ZipFile(weighs, 'r') as zip_ref:
11     zip_ref.extractall(extract_path)
```

Extracto de código 5.2: Código para importar las imágenes

Una vez teníamos ese dataset en Colab ya podemos hacer uso de nuestras imágenes pero antes tuvimos que preprocesarlas para reducir la dimensión de 512x512x3 a 64x64x3 para poder usarlas como inputs en el autoencoder.

Esto lo hicimos mediante el siguiente código 5.3:

```
1 # Define el path a la carpeta donde estn las imagenes
2 training_autoencoder_folder_path = '/content/sl_frames/training_set'
3
4 # Obtiene una lista de los nombres de las imagenes
5 image_files = os.listdir(training_autoencoder_folder_path)
6
7 # Carga las imagenes en una lista
8 image_list = []
9 for filename in image_files:
10     img = Image.open(os.path.join(training_autoencoder_folder_path, filename))
11     img = img.resize((64, 64))
12     img = np.array(img)
13     image_list.append(img)
14
15 # Convierte la lista en un array que puede procesar el autoencoder
16 X_train_autoencoder = np.array(image_list)
17
18 print(X_train_autoencoder.shape)
```

Extracto de código 5.3: Resize de las imágenes

Una vez detectamos el problema de que estábamos realizando los experimentos con un dataset que no correspondía, cambiamos de dataset al propio y usamos la estructura de importe ya montada para usarlos.

5.4.2. Versiones y entrenamiento de los modelos

Vamos a pasar ahora a lo que más tiempo nos ha llevado del trabajo, generar versiones de los autoencoder que funcionasen y además que consiguiesen entrenar hasta minimizar la función de pérdida hasta valores aceptables.

Versión 1

La primera versión del autoencoder que propusimos fue la siguiente: 5.4. Esta versión añade capas de pooling y upsample a la estructura básica del autoencoder propuesta en [1].

Los problemas que presentaba esta versión es que al tener un mayor número de capas los entrenamientos eran muy largos llegando a demorar hasta cuatro horas ya que para poder bajar la pérdida en la reconstrucción tuvimos que entrenar por 300 epochs consiguiendo bajarla hasta 17 tanto en el conjunto de entrenamiento como en el de validación.

```

1 # Encoder
2 encoder_input = tf.keras.Input(shape=(64, 64, 3))
3
4 x = layers.Conv2D(128, (3, 3), activation=tf.nn.leaky_relu, padding='same', name="conv1")(encoder_input)
5 x = MaxPooling2D((2, 2), padding='same', name="pool1")(x)
6
7 x = layers.Conv2D(256, (3, 3), activation=tf.nn.leaky_relu, padding='same', name="conv2")(x)
8 x = MaxPooling2D((2, 2), padding='same', name="pool2")(x)
9
10 x = layers.Conv2D(512, (3, 3), activation=tf.nn.leaky_relu, padding='same', name="conv3")(x)
11 x = MaxPooling2D((2, 2), padding='same', name="pool3")(x)
12
13 x = layers.Conv2D(1024, (3, 3), activation=tf.nn.leaky_relu, padding='same', name="conv4")(x)
14 x = MaxPooling2D((2, 2), padding='same', name="pool4")(x)
15
16 x = layers.Flatten()(x)
17 encoded = layers.Dense(128, activation="linear", name="dense_layer")(x)
18 encoder = Model(encoder_input, encoded, name='encoder')
19
20
21 # Decoder
22 decoder_input = layers.Input(shape=(128,))
23
24 x = layers.Dense(16384, activation=tf.nn.leaky_relu)(decoder_input)
25 x = layers.Reshape((8, 8, 256))(x)
26
27 x = layers.Conv2DTranspose(512, (3, 3), padding='same', activation=tf.nn.leaky_relu, name="conv5")(x)
28 x = UpSampling2D((2, 2), name="up1")(x)
29
30 x = layers.Conv2DTranspose(256, (3, 3), padding='same', activation=tf.nn.leaky_relu, name="conv6")(x)
31 x = UpSampling2D((2, 2), name="up2")(x)
32
33 x = layers.Conv2DTranspose(128, (3, 3), padding='same', activation=tf.nn.leaky_relu, name="conv7")(x)
34 x = UpSampling2D((2, 2), name="up3")(x)
35
36 decoded = layers.Conv2DTranspose(3, (3, 3), activation='linear', padding='same', name="final_conv")(x)
37
38
39 decoder = Model(decoder_input, decoded, name='decoder')
40
41 # Autoencoder model
42 autoencoder_input = tf.keras.Input(shape=(64, 64, 3))
43 encoded_img = encoder(autoencoder_input)
44 decoded_img = decoder(encoded_img)
45 autoencoder = Model(autoencoder_input, decoded_img, name='autoencoder')
46
47 opt = keras.optimizers.Adam(learning_rate=0.001)
48 autoencoder.compile(optimizer=opt, loss='mse') # Using mean squared error as the loss function

```

Extracto de código 5.4: Autoencoderv1

Sin embargo a pesar de que conseguía clasificar el conjunto de testeo con un 65 por ciento de precisión, es decir clasificaba las imágenes no anómalas como 0 en la mayoría de los casos y las anómalas como 1. Descartamos este modelo ya que al usar la métrica AUC [1, Página 6, Apartado F] obteníamos un valor de 0.5. La métrica AUC es una medida utilizada para evaluar el rendimiento de un modelo de clasificación binaria. Esta muestra un gráfico que enfrenta los verdaderos positivos frente a los falsos positivos. Obtener una puntuación de un 1 sería tener un clasificador perfecto y un 0.5 representaría acertar solo la mitad de las veces.

Versión 2

La siguiente versión es la traducción directa capa por capa del modelo 5.5. Esta traduce directamente las capas del modelo original a Tensorflow.

```
1 # ENCODER
2 encoder = keras.Sequential(name='encoder')
3
4 encoder.add(keras.layers.Conv2D(filters=128, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu", name
   ='input', input_shape=(64, 64, 3)))
5
6 encoder.add(keras.layers.Conv2D(filters=256, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu", name
   ='conv1'))
7
8 encoder.add(keras.layers.Conv2D(filters=512, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu", name
   ='conv2'))
9
10 encoder.add(keras.layers.Conv2D(filters=1024, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu",
   name='conv3'))
11
12 encoder.add(keras.layers.Flatten(name='flatten'))
13 encoder.add(keras.layers.Dense(units=128, activation='linear', name='dense'))
14 encoder.add(keras.layers.Embedding(input_dim=128, output_dim=128, embeddings_initializer='glorot_uniform',
   embeddings_regularizer='L1L2'))
15
16 # DECODER
17 decoder = keras.Sequential(name='decoder')
18
19 decoder.add(keras.layers.Input(shape=(128, 128), name='input'))
20 decoder.add(keras.layers.Reshape((4, 4, 1024)))
21
22 decoder.add(keras.layers.Conv2D(filters=512, kernel_size=(3, 3), activation='leaky_relu', padding='same'))
23
24 decoder.add(keras.layers.Conv2DTranspose(filters=256, kernel_size=(3, 3), strides=(2, 2), activation="
   leaky_relu", padding='same', name='deconvconv1'))
25
26 decoder.add(keras.layers.Conv2DTranspose(filters=128, kernel_size=(3, 3), strides=(2, 2), activation="
   leaky_relu", padding='same', name='deconvconv2'))
27
28 decoder.add(keras.layers.Conv2DTranspose(filters=64, kernel_size=(3, 3), strides=(2, 2), activation="
   leaky_relu", padding='same', name='deconvconv3'))
29
30 decoder.add(keras.layers.Conv2DTranspose(filters=3, kernel_size=(3, 3), strides=(2, 2), activation="linear",
   padding='same', name='output'))
31
32 # AUTOENCODER
33 autoencoder_input = keras.layers.Input(shape=(64, 64, 3), name='autoencoder_input')
34 encoded_repr = encoder(autoencoder_input)
35 decoded_repr = decoder(encoded_repr)
36
37 autoencoder = keras.models.Model(inputs=autoencoder_input, outputs=decoded_repr, name='autoencoder')
38
39 opt = keras.optimizers.Adam(learning_rate=0.001)
40 autoencoder.compile(optimizer=opt, loss='mse')
```

Extracto de código 5.5: Autoencoder v2

El problema que presentó esta arquitectura fue que a partir de tempranos epochs el modelo se sobreajustaba demasiado al conjunto de entrenamiento. Cuando el entrenamiento llevaba aproximadamente 15 epochs dejaba de entrenar y la pérdida se estancaba en 50 aproximadamente. Por ello pasamos a la tercera versión.

Versión 3

La tercera versión ?? implementa la regularización L1L2 en la capa en la que se genera la codificación. Simultáneamente se implementa en la misma capa la inicialización de los pesos de esta. Inicialmente se probó inicializarlos con *he-normal* y conseguimos de esta manera retrasar el overfitting hasta el epoch 100 sin embargo la pérdida quedaba en 45 que seguía siendo elevada para los resultados que queremos obtener. Luego probamos la inicialización de los pesos con *glorot-uniform* y reertasamos el overfitting a los 115 epochs. Sin embargo, la pérdida se estancaba aproximadamente en el mismo valor.

```
1 # ENCODER
2 encoder = keras.Sequential(name='encoder')
3
4 encoder.add(keras.layers.Conv2D(filters=128, kernel_size=(3, 3), strides=(2, 2),
5                                 activation="leaky_relu", name='input', input_shape=(64, 64, 3)))
6
7 encoder.add(keras.layers.Conv2D(filters=256, kernel_size=(3, 3), strides=(2, 2),
8                                 activation="leaky_relu", name='conv1'))
9
10 encoder.add(keras.layers.Conv2D(filters=512, kernel_size=(3, 3), strides=(2, 2),
11                                activation="leaky_relu", name='conv2'))
12
13 encoder.add(keras.layers.Conv2D(filters=1024, kernel_size=(3, 3), strides=(2, 2),
14                                activation="leaky_relu", name='conv3',
15                                kernel_regularizer=regularizers.l2(0.01)))
16
17 encoder.add(keras.layers.Flatten(name='flatten'))
18 encoder.add(keras.layers.Dense(units=128, activation='linear', name='dense',
19                                kernel_regularizer=regularizers.l2(0.01)))
20 encoder.add(keras.layers.Embedding(input_dim=128, output_dim=128, embeddings_initializer='glorot_uniform',
21                                    embeddings_regularizer='L1L2'))
22
23 # DECODER
24 decoder = keras.Sequential(name='decoder')
25
26 decoder.add(keras.layers.Input(shape=(128, 128), name='input'))
27 decoder.add(keras.layers.Reshape((4, 4, 1024)))
28
29 decoder.add(keras.layers.Conv2D(filters=512, kernel_size=(3, 3), activation='leaky_relu', padding='same'))
30
31 decoder.add(keras.layers.Conv2DTranspose(filters=256, kernel_size=(3, 3), strides=(2, 2),
32                                           activation="leaky_relu", padding='same', name='deconvconv1',
33                                           kernel_regularizer=regularizers.l2(0.01)))
34
35 decoder.add(keras.layers.Conv2DTranspose(filters=128, kernel_size=(3, 3), strides=(2, 2),
36                                           activation="leaky_relu", padding='same', name='deconvconv2'))
37
38 decoder.add(keras.layers.Conv2DTranspose(filters=64, kernel_size=(3, 3), strides=(2, 2),
39                                           activation="leaky_relu", padding='same', name='deconvconv3'))
40
41 decoder.add(keras.layers.Conv2DTranspose(filters=3, kernel_size=(3, 3), strides=(2, 2),
42                                           activation="linear", padding='same', name='output'))
43
44 # AUTOENCODER
45 autoencoder_input = keras.layers.Input(shape=(64, 64, 3), name='autoencoder_input')
46 encoded_repr = encoder(autoencoder_input)
47 decoded_repr = decoder(encoded_repr)
48 autoencoder = keras.models.Model(inputs=autoencoder_input, outputs=decoded_repr, name='autoencoder')
49 opt = keras.optimizers.Adam(learning_rate=0.001)
50 autoencoder.compile(optimizer=opt, loss='mse')
```

Extracto de código 5.6: Autoencoder v3

También se probó a aplicar dropout y data augmentation. Para ambas aplicaciones el modelo no era capaz de reducir la función de pérdida.

El aumento de datos se hizo de la siguiente manera [5.7](#). Aplicamos rotaciones de 40°, flip horizontal, cambios al brillo y zoom:

```
1 images = X_train_autoencoder
2 data_format = 'channels_last'
3
4 datagen = ImageDataGenerator(
5     rotation_range=40,
6     horizontal_flip=True,
7     brightness_range=[0.6, 1.7],
8     zoom_range=[0.5, 1.0],
9     fill_mode='nearest'
10 )
11
12 augmented_images = datagen.flow(images, batch_size=len(images), shuffle=False).next()
```

Extracto de código 5.7: Data augmentation

Aun así debido a los pobres resultados tuvimos que tomar otra alternativa que fue probar a entrenar el modelo mediante la técnica de validación cruzada [5.8](#).

```
1 # Parmetros de cross-validation
2 n_splits = 10
3 kf = KFold(n_splits=n_splits, shuffle=True, random_state=42)
4
5 model = autoencoder
6
7 mse_scores = []
8 i = 0
9 for train_index, val_index in kf.split(X_train_autoencoder):
10     early_stopping = EarlyStopping(monitor='val_loss', patience=5, restore_best_weights=True)
11     reduce_lr = ReduceLROnPlateau(monitor='val_loss', factor=0.1, patience=4, min_lr=1e-7)
12
13     X_train_fold, X_val_fold = X_train_autoencoder[train_index], X_train_autoencoder[val_index]
14     # Entrenar el modelo
15     model.fit(X_train_fold, X_train_fold, epochs=10, batch_size=32, validation_data=(X_val_fold, X_val_fold),
16             verbose=1, callbacks=[early_stopping, reduce_lr])
17
18     # Evaluar el modelo
19     mse_loss = model.evaluate(X_val_fold, X_val_fold, verbose=0)
20     mse_scores.append(mse_loss)
21     i += 1
22     print(f'Vamos por el split {i}')
23 print(f"Mean MSE over {n_splits} folds: {np.mean(mse_scores)}")
```

Extracto de código 5.8: Cross Validation

Esta técnica consiguió aumentar los epochs en los que se producía el overfitting considerablemente. Sin embargo el modelo seguía estancado en la misma pérdida. Por tanto pasamos a la siguiente versión.

Versión 4

La cuarta versión del autoencoder^{5.9} elimina la capa de embeddings y genera la codificación directamente de un vector 128 como hacíamos en la versión original. De esta forma y aplicando la inicialización de pesos en *he-normal* y *gradient-clipping*, sumado al uso de *cross-validation* para el entrenamiento consiguieron bajar la pérdida hasta el valor de 35. De todas formas, aún se podía bajar aún más ya que con la primera versión conseguimos llegar a 17.

```
1 # Normalización de la entrada (suponiendo que tus datos están en formato numpy array)
2
3 # Supongamos que X_train es tu conjunto de datos de entrenamiento
4
5 # ENCODER
6 encoder = models.Sequential(name='encoder')
7 encoder.add(layers.Conv2D(filters=64, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu",
8     kernel_initializer='glorot_uniform', kernel_regularizer=regularizers.l2(0.01), input_shape=(64, 64, 3),
9     name='input'))
10 encoder.add(layers.Conv2D(filters=128, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu",
11     kernel_initializer='glorot_uniform', kernel_regularizer=regularizers.l2(0.01), name='conv1'))
12 encoder.add(layers.Conv2D(filters=256, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu",
13     kernel_initializer='glorot_uniform', kernel_regularizer=regularizers.l2(0.01), name='conv2'))
14 encoder.add(layers.Flatten(name='flatten'))
15 encoder.add(layers.Dense(units=128, activation='linear', kernel_initializer='glorot_uniform',
16     kernel_regularizer=regularizers.l2(0.01), name='dense'))
17
18 # DECODER
19 decoder = models.Sequential(name='decoder')
20 decoder.add(layers.Input(shape=(128,), name='input'))
21 decoder.add(layers.Dense(4*4*256, activation='linear', kernel_initializer='glorot_uniform', kernel_regularizer=
22     regularizers.l2(0.01), name='dense'))
23 decoder.add(layers.Reshape((4, 4, 256)))
24 decoder.add(layers.Conv2DTranspose(filters=128, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu",
25     kernel_initializer='glorot_uniform', kernel_regularizer=regularizers.l2(0.01), padding='same', name='
26     deconvconv1'))
27 decoder.add(layers.Conv2DTranspose(filters=64, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu",
28     kernel_initializer='glorot_uniform', kernel_regularizer=regularizers.l2(0.01), padding='same', name='
29     deconvconv2'))
30 decoder.add(layers.Conv2DTranspose(filters=32, kernel_size=(3, 3), strides=(2, 2), activation="leaky_relu",
31     kernel_initializer='glorot_uniform', kernel_regularizer=regularizers.l2(0.01), padding='same', name='
32     deconvconv3'))
33 decoder.add(layers.Conv2DTranspose(filters=3, kernel_size=(3, 3), strides=(2, 2), activation="linear",
34     kernel_initializer='glorot_uniform', kernel_regularizer=regularizers.l2(0.01), padding='same', name='
35     output'))
36
37 # AUTOENCODER
38 autoencoder_input = layers.Input(shape=(64, 64, 3), name='autoencoder_input')
39 encoded_repr = encoder(autoencoder_input)
40 decoded_repr = decoder(encoded_repr)
41
42 autoencoder = models.Model(inputs=autoencoder_input, outputs=decoded_repr, name='autoencoder')
43
44 # Compile the model
45 opt = Adam(learning_rate=0.001, clipnorm=1.0)
46 autoencoder.compile(optimizer=opt, loss='mse')
```

Extracto de código 5.9: Autoencoder v4

Versión 5

Pasamos así a la versión final del autoencoder ?? En esta versión eliminamos capas de la estructura del autoencoder para así reducir los parametros de la red y reducir el sobreajuste. Esto conbinado con la inicialización de los pesos de las capas densas con *glorot-uniform* y el uso de *cross-validation* en el entrenamiento consiguió bajar la pérdida del modelo a 22. Decidimos aceptarlo ya que probando con las diferentes combinatorias durante los entrenamientos no conseguimos obtener mejores resultados.

```
1 encoder = keras.Sequential(name='encoder')
2
3 encoder.add(keras.layers.Conv2D(filters=32, kernel_size=(3, 3),strides=(2, 2),padding='same',activation='
    leaky_relu',name='conv1',input_shape=(64, 64, 3)))
4
5 encoder.add(keras.layers.Conv2D(filters=64,kernel_size=(3, 3),strides=(2, 2),padding='same',activation='
    leaky_relu',name='conv2',
6 ))
7
8 encoder.add(keras.layers.Conv2D(filters=128,kernel_size=(3, 3),strides=(2, 2),padding='same',activation='
    leaky_relu',name='conv3',
9 ))
10
11 encoder.add(layers.Flatten(name='flatten'))
12 encoder.add(layers.Dense(units=128, activation='linear', kernel_initializer='glorot_uniform', name='dense'))
13
14 decoder = keras.Sequential(name='decoder')
15 decoder.add(layers.Dense(units=8 * 8 * 128, activation='linear', kernel_initializer='glorot_uniform', name='
    dense2'))
16 decoder.add(layers.Reshape((8, 8, 128), name='reshape'))
17
18 decoder.add(keras.layers.Conv2DTranspose(filters=64, kernel_size=(3, 3),strides=(2, 2), padding='same',
    activation='leaky_relu', name='conv5',
19 ))
20
21 decoder.add(keras.layers.Conv2DTranspose( filters=32, kernel_size=(3, 3),strides=(2, 2),padding='same',
    activation='leaky_relu',name='conv6',
22 ))
23
24 decoder.add(keras.layers.Conv2DTranspose(filters=3,kernel_size=(3, 3),strides=(2, 2),padding='same',activation
    ='linear',name='output'
25 ))
26
27 autoencoder_input = layers.Input(shape=(64, 64, 3), name='autoencoder_input')
28 encoded_repr = encoder(autoencoder_input)
29 decoded_repr = decoder(encoded_repr)
30
31 autoencoder = models.Model(inputs=autoencoder_input, outputs=decoded_repr, name='autoencoder')
32 opt = Adam(learning_rate=0.001, clipnorm=1.0)
33 autoencoder.compile(optimizer=opt, loss='mse')
```

Extracto de código 5.10: Código para importar las imágenes

El entrenamiento de las versiones 2 en adelante tarda aproximadamene unos 25 minutos con el entorno L4, que es de pago en Colab. Esto mejora considerablemente las 4 horas de entrenamiento de la versión 1. El motivo de la adquisición de GoogleColab premium fue que, en la versión gratuita, la gráfica se quedaba sin memoria al ser modelos de un tamaño relativamente elevado.

5.4.3. Ejemplo de uso del autoencoder

Como hemos establecido, queremos utilizar el autoencoder como detector de anomalías utilizando el error en la reconstrucción de la imagen a modo de puntuación de anomalía. Si la imagen no tiene anomalía la reconstrucción cometerá un error bajo. Por otro lado, si la imagen contiene anomalías la reconstrucción cometerá un error elevado ya que no reconstruirá la anomalía en sí debido a que no está entrenado para ello.

Vamos a ver ahora como realiza estas reconstrucciones [5.11](#):

```
1 image_array1 = img
2 image_array2 = test1_decode[0]
3
4 fig, axes = plt.subplots(1, 2)
5
6 axes[0].imshow(image_array1.astype('uint8'))
7 axes[0].axis('off')
8 axes[0].set_title('Original')
9 axes[1].imshow(image_array2.astype('uint8'))
10 axes[1].axis('off')
11 axes[1].set_title('Reconstruida')
12 plt.show()
```

Extracto de código 5.11: Reconstrucción sin anomalía

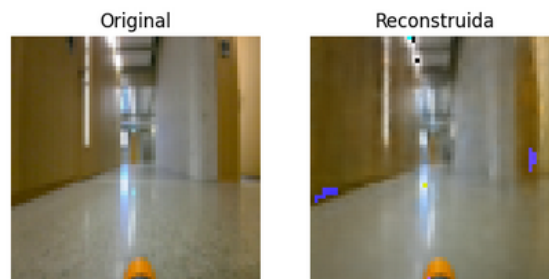


Figura 5.3: Reconstrucción sin anomalía

Como podemos ver la reconstrucción de la imagen es bastante fiable. Concretamente el error cometido es de un valor de 46.

Veamos ahora la reconstrucción de una imagen que contiene una anomalía [5.12](#).

```
1 image_array1 = img2
2 image_array2 = test2_decode[0]
3
4 fig, axes = plt.subplots(1, 2)
5
6 axes[0].imshow(image_array1.astype('uint8'))
7 axes[0].axis('off')
8 axes[0].set_title('Original')
9
10 axes[1].imshow(image_array2.astype('uint8'))
11 axes[1].axis('off')
12 axes[1].set_title('Reconstruida')
13
14 plt.show()
```

Extracto de código 5.12: Reconstrucción con anomalía

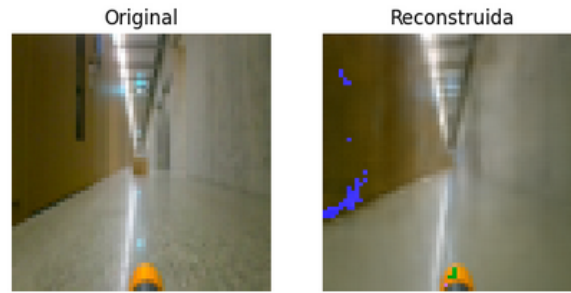


Figura 5.4: Reconstrucción con anomalía

Como podemos ver al reconstruir la imagen la caja que se observa en la imagen original ha sido eliminada en la reconstrucción haciendo que el error cometido aumente considerablemente. Este aumento es lo que usaremos como indicador de que en la imagen hay una anomalía.

5.4.4. Autoencoder como detector de anomalías

Para poder usar el autoencoder de esta manera hemos creado un par de funciones auxiliares que dada una imagen, predice si es una anomalía.

La primera función simplemente calcula el error medio cuadrático que se produce al reconstruir la imagen (que es la función que minimiza la pérdida).

```

1 def mse_autoencoder_predict(image, path):
2     img = Image.open(os.path.join(path, image))
3     img = img.resize((64, 64))
4     img = np.array(img)
5     img_array_reshaped = np.expand_dims(img, axis=0)
6     encoder_model = autoencoder.get_layer('encoder')
7     decoder_model = autoencoder.get_layer('decoder')
8     encoded_img = encoder_model.predict(img_array_reshaped)
9     decoded_img = decoder_model.predict(encoded_img)
10    mse = np.mean(np.square(img - decoded_img[0]))
11    return mse

```

Extracto de código 5.13: Error cuadrático medio

La función codifica, decodifica y finalmente calcula el error cometido. Vamos a mostrar ahora dos ejemplos 5.5. En el primero la imagen contiene una anomalía y en el segundo no. Observamos como el error del primer cálculo es mucho mayor.

```
[ ] mse_autoencoder_predict('010775_512_512.jpg',test_models_folder_path)
```

```
1/1 [=====] - 0s 22ms/step  
1/1 [=====] - 0s 20ms/step  
232.28996
```

```
[ ] mse_autoencoder_predict('000700_512_512.jpg',test_models_folder_path)
```

```
1/1 [=====] - 0s 23ms/step  
1/1 [=====] - 0s 21ms/step  
60.518055
```

Figura 5.5: Ejemplos error cuadrático medio

La siguiente función 5.14 es la que realiza la clasificación de la imagen. Simplemente si la función 5.13 devuelve un error mayor a un umbral, clasifica la imagen como anomalía. En este caso a base de diferentes testeos hemos establecido que el umbral de detección sea 200 ya que es con el que hemos obtenido los mejores resultados.

```
1 def classifier_AE(name_list,path):  
2     y_predict = []  
3     for i in range(len(name_list)):  
4         mse = mse_autoencoder_predict(name_list[i],path)  
5         print(mse)  
6         if mse >= 200:  
7             y_predict.append(1)  
8         else:  
9             y_predict.append(0)  
10        print(i)  
11    return y_predict
```

Extracto de código 5.14: Clasificador autoencoder

Para predecir el conjunto de datos entero con un bucle for que computase ambas funciones tardamos 8 horas. Por lo tanto hicimos una versión de las funciones de forma que hiciera la predicción del conjunto de datos por batches 5.15. De esta forma la predicción pasó a tardar 2 horas.

```

1
2 def mse_autoencoder_predict(images, path, autoencoder):
3     imgs = []
4     for image in images:
5         img = Image.open(os.path.join(path, image))
6         img = img.resize((64, 64))
7         img = np.array(img)
8         imgs.append(img)
9
10    imgs_array = np.array(imgs)
11    encoder_model = autoencoder.get_layer('encoder')
12    decoder_model = autoencoder.get_layer('decoder')
13    encoded_imgs = encoder_model.predict(imgs_array)
14    decoded_imgs = decoder_model.predict(encoded_imgs)
15
16    mses = []
17    for img, decoded_img in zip(imgs_array, decoded_imgs):
18        img_flat = img.flatten()
19        decoded_img_flat = decoded_img.flatten()
20        mse = mean_squared_error(img_flat, decoded_img_flat)
21        mses.append(mse)
22
23    return mses
24
25 def classifier_AE(name_list, path, autoencoder, batch_size=32):
26     y_predict = []
27     for i in range(0, len(name_list), batch_size):
28         batch_names = name_list[i:i+batch_size]
29         mses = mse_autoencoder_predict(batch_names, path, autoencoder)
30         print(batch_names)
31         print(mses)
32         for mse in mses:
33             if mse >= 200:
34                 y_predict.append(1)
35             else:
36                 y_predict.append(0)
37         print(f'Processed batch {i//batch_size + 1}')
38     return y_predict

```

Extracto de código 5.15: Clasificador autoencoder por batches

Una vez hechas las predicciones calculamos el porcentaje de acierto y obtuvimos aproximadamente un 63 por ciento, valor muy parecido a la primera versión.

Desafortunadamente al aplicar de nuevo la métrica de AUC, volvemos a obtener un 0.5.

Esto es debido a que a pesar de que clasifica correctamente de forma general, cuando las anomalías son muy pequeñas aun en las imágenes no es capaz de detectar que están ahí ya que la reconstrucción es prácticamente igual y no sobrepasa el umbral en la clasificación.

5.4.5. Codificación de las imágenes

Para realizar la codificación de las imágenes del conjunto de testeo simplemente aplicamos un bucle for que va codificando por batches ??.

```
1 images = X_train_models
2
3 encoder_model = autoencoder.get_layer('encoder')
4
5 batch_size = 32
6 encoded_images1 = []
7
8 for i in range(0, len(images), batch_size):
9     batch_images1 = images[i:i+batch_size]
10    encoded_imgs1 = encoder_model.predict(batch_images1)
11    encoded_images1.append(encoded_imgs1)
12
13 encoded_images1 = np.concatenate(encoded_images1, axis=0)
14
15
16 print(encoded_images1.shape)
```

Extracto de código 5.16: Codificación de las imágenes

Hacemos esto tanto sobre el conjunto de entrenamiento del perceptrón multicapa como sobre el conjunto de testeo para realizar las pruebas finales.

5.4.6. Perceptrón multicapa

Como vimos en 4.3 la estructura del perceptrón es muy simple constando solo de tres capas densas que hará de clasificador de las imágenes codificadas.

El perceptrón obtiene una reducción de su función de error de 2.3941e-12 obteniendo un accuracy del 100 por cien.

Sin embargo al evaluar el modelo bajo el AUC, volvemos a obtener un 0.5 debido a que acarreamos el problema del autoencoder en la codificación. Los vectores generados en los que las anomalías son muy pequeñas acaban clasificados como no anómalos en vez de anómalos.

6. Conclusiones

Tras analizar la arquitectura propuesta en [1] tratamos de realizar una traducción de su código que se encuentra en el siguiente [github](#), hemos conseguido producir dos modelos de predicción de anomalías que, cuando estas se hacen realmente notables en la imagen son capaces de detectarlas. En primer lugar, hemos implementado varias versiones del autoencoder aplicando las diferentes técnicas mencionadas para reducir el sobreajuste de los entrenamientos, consiguiendo mejorar la primera versión tanto en tiempo de entrenamiento como en su capacidad de reconstrucción de las imágenes. Por otra parte hemos combinado esta primera arquitectura con un perceptrón multicapa, dando así lugar a una segunda arquitectura con la misma funcionalidad y un rendimiento similar.

Durante el desarrollo de este trabajo se ha hecho una revisión bibliográfica que nos ha permitido llevar a cabo la implementación y entrenamiento de los modelos presentados aplicando numerosas técnicas que se usan en deep learning para reducir el sobreajuste de estos. Algunas de estas técnicas ya se usaban en la versión del pape, como el data augmentation. Sin embargo muchas otras mencionadas en la sección 4 han sido buscadas e implementadas desde cero.

Además, hemos montado un entorno explotando al máximo las funcionalidades de GoogleDrive y GoogleColab capaz de soportar las necesidades de este proyecto, proceso que ha presentado numerosas complicaciones que hemos podido ir solventando.

El notebook en el que se realizan los experimentos queda publicado en el siguiente [github](#).

Bibliografía

- [1] Dario Mantegazza, Alessandro Giusti, Luca Maria Gambardella, and Jérôme Guzzi. An outlier exposure approach to improve visual anomaly detection performance for mobile robots. *IEEE Robotics and Automation Letters*, 7(4): 11354–11361, 2022.
- [2] PyTorch. Pytorch, 2024. URL <https://pytorch.org/>.
- [3] TensorFlow. Tensorflow, 2024. URL <https://www.tensorflow.org/>.
- [4] Google Colab. Google colaboratory, 2024. URL <https://colab.research.google.com/>.
- [5] Talent.com. Sueldo medio de un ingeniero informático, 2024. URL <https://es.talent.com/salary?job=ingeniero+inform%C3%A1tico>.
- [6] Francois Chollet. *Deep learning with Python*. Simon and Schuster, 2021.
- [7] Yifei Zhang. A better autoencoder for image: Convolutional autoencoder. In *ICONIP17-DCEC*. Available online: http://users.cecs.anu.edu.au/Tom.Gedeon/conf/ABCS2018/paper/ABCS2018_paper_58.pdf (accessed on 23 March 2017), 2018.
- [8] Michael Cogswell, Faruk Ahmed, Ross Girshick, Larry Zitnick, and Dhruv Batra. Reducing overfitting in deep networks by decorrelating representations. *arXiv preprint arXiv:1511.06068*, 2015.
- [9] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of big data*, 6(1):1–48, 2019.
- [10] Mustafa Abdul Salam, Ahmad Taher Azar, Mustafa Samy Elgendy, and Khaled Mohamed Fouad. The effect of different dimensionality reduction techniques on machine learning overfitting problem. *Int. J. Adv. Comput. Sci. Appl*, 12(4):641–655, 2021.