



**COMILLAS**  
UNIVERSIDAD PONTIFICIA

ICAI

**GRADO EN INGENIERÍA  
MATEMÁTICA E INTELIGENCIA  
ARTIFICIAL**

**TRABAJO FIN DE GRADO**

Generative AI for Medical Images: Evaluating  
Synthetic Data in Diagnostic Model

Autor: Ignacio Viadero Canduela

Director: Sergio Altares López

Co-Director: Jose María Bengochea Guevara

Madrid, Junio 2026



Declaro, bajo mi responsabilidad, que el Proyecto presentado con el título  
en la ETS de Ingeniería - ICAI de la Universidad Pontificia Comillas en el  
curso académico 2025/ es de mi autoría, original e inédito y  
no ha sido presentado con anterioridad a otros efectos.  
El Proyecto no es plagio de otro, ni total ni parcialmente, y la información que  
ha sido tomada de otros documentos está debidamente referenciada.



Fdo.: Ignacio Viadero Canduela

Fecha: 15 / 06 / 2026

Autorizada la entrega del proyecto  
EL DIRECTOR DEL PROYECTO

Fdo.: Sergio Altares López

Fecha: 15 / 06 / 2026

EL CODIRECTOR DEL PROYECTO

Fdo.: Jose María Bengochea Guevara

Fecha: 15 / 06 / 2026

# Agradecimientos

No se puede empezar de otra manera que agradeciendo a toda la gente que ha permitido la realización y el logro de este proyecto, y que me ha ayudado a lo largo de mi camino hasta aquí. Con este trabajo se cierra una etapa de mi vida, llena de momentos gratificantes y algún que otro altibajo, pero como dijo Alexander Bell: “Cuando una puerta se cierra, otra se abre”, y espero ansioso a la llegada de nuevas experiencias y etapas en las que aprenderé y mejoraré tanto profesional como personalmente como he hecho hasta ahora.

En primer lugar, me gustaría agradecer a mis dos tutores del Trabajo de Fin de Grado, José María Bengochea Guevara y Sergio Altares López. A lo largo de todo el proyecto han estado disponibles siempre que lo he necesitado, resolviendo dudas, aportando ideas y ayudándome a tomar decisiones cuando el camino no estaba claro. Sin su ayuda, paciencia y dedicación este trabajo no habría sido posible, y les agradezco especialmente la confianza que han depositado en mí y en este proyecto desde el principio. Cabe mencionar también a todos los profesores y a mi Universidad Pontificia Comillas, cuya enseñanza ha sido fundamental para poder llevar a cabo este trabajo.

Por otro lado, me gustaría mencionar a mis compañeros y amigos de clase, han estado al lado mío a lo largo de estos 4 años de carrera, donde hemos celebrado juntos nuestras victorias y nos hemos ayudado para superar nuestras mayores dificultades. Me han acompañado durante todas esas horas de estudio, en las épocas de exámenes e incluso en todos los proyectos que tengo la suerte de haber podido compartir con ellos, por eso creo que una parte de este trabajo es gracias a ellos, y espero que sigan a mi lado en los nuevos retos que nos esperan.

Por último, me gustaría agradecer a mi familia y a mi novia, que han estado ahí siempre que he necesitado a alguien que me tendiera la mano para levantarme y seguir adelante. Han sido mi apoyo en los momentos más difíciles y también los primeros en celebrar cada uno de mis logros, y por ello este trabajo también es un poco suyo. A mi novia por ayudarme y estar a mi lado durante todos estos momentos, por sacarme una sonrisa siempre que la necesitaba y guiarme a conseguir mi mejor versión. Gracias por la paciencia en los días complicados, por escucharme incluso cuando ni yo mismo sabía

explicarme y por el apoyo incondicional. A mis hermanos, que han sido una motivación extra para seguir adelante y dar lo mejor de mí. Han sido un ejemplo y un impulso para sacar lo mejor de mi, compartir este camino con ellos ha hecho que todo sea mucho más llevadero y significativo. Y por último, a mis padres, gracias por hacer el esfuerzo emocional y económico de darme la oportunidad de acudir a esta grandísima universidad y de vivir la experiencia universitaria en Madrid, que han depositado siempre toda su confianza en mí y no dudaron que este era mi camino a seguir. Gracias por creer siempre en mí, por vuestro apoyo constante y por enseñarme el valor del esfuerzo y la constancia, que han sido fundamentales para llegar hasta aquí.

Quiero dedicar también este trabajo a todas las personas que conviven con la enfermedad de Alzheimer y, en especial, a los pacientes, víctimas, cuidadores y familiares que sufren de cerca su impacto cada día. Detrás de cada avance científico, de cada experimento y de cada resultado, hay una realidad profundamente humana marcada por el esfuerzo, la incertidumbre, el desgaste emocional y, muchas veces, el dolor. Este proyecto nace también con la esperanza de que, desde la ingeniería y la investigación, se pueda aportar aunque sea un pequeño paso hacia una mejor comprensión, un diagnóstico más temprano y, en el futuro, una mejor calidad de vida para quienes se enfrentan a esta enfermedad. A todos ellos, por su lucha, su fortaleza y su ejemplo, va dedicada una parte muy especial de este trabajo.

# Resumen

La enfermedad del Alzheimer es una patología neurodegenerativa que supone uno de los principales retos de la neurología actual, ya que a pesar de décadas de investigación, los mecanismos que causan la enfermedad de Alzheimer no se conocen por completo, y siguen existiendo importantes interrogantes sobre su origen y progresión. Una de las razones que limita el avance en este campo es la escasez y el desequilibrio de los conjuntos de datos disponibles, especialmente en fases muy tempranas de la enfermedad, donde la detección de patrones es difícil y la recogida de información se ve restringida además por las políticas de protección de datos y por la propia vulnerabilidad de los pacientes. Esta variabilidad se ve incrementada por la influencia de múltiples factores internos y externos del individuo y por la evolución asimétrica, de modo que la enfermedad no sigue siempre un patrón único y reproducible, lo que complica aún más la identificación de biomarcadores robustos y generalizables.

Este estudio tiene como objetivo desarrollar y evaluar un sistema basado en un modelo de aprendizaje profundo con dos fases claramente diferenciadas. En una primera fase se entrena un modelo de clasificación capaz de clasificar cada sujeto a su estado cognitivo a partir de resonancias magnéticas (MRI). En la segunda fase se entrena un modelo generativo condicionado por clase, encargado de sintetizar nuevas imágenes realistas a partir del mismo conjunto finito de volúmenes, y se analiza hasta qué punto la incorporación de estas imágenes artificiales al conjunto de entrenamiento permite mejorar el rendimiento del clasificador original y mitigar, al menos parcialmente, el problema de la escasez de datos en el ámbito médico. Finalmente, se incluye una componente de explicabilidad que permite interpretar tanto las decisiones del modelo de clasificación como las diferencias morfológicas entre clases, facilitando así una lectura más clínica de los resultados obtenidos.

# Índice

<b>1</b>	<b>Introducción y motivación</b>	<b>1</b>
1.1	Contexto general del Alzheimer y diagnóstico por imagen . . . . .	1
1.2	Motivación . . . . .	2
1.3	Planteamiento del problema y preguntas de investigación . . . . .	2
<b>2</b>	<b>Objetivos y contribuciones</b>	<b>4</b>
2.1	Objetivo general . . . . .	4
2.2	Objetivos específicos . . . . .	4
2.3	Contribuciones . . . . .	5
<b>3</b>	<b>Fundamentos teóricos</b>	<b>6</b>
3.1	Enfermedad del Alzheimer . . . . .	6
3.1.1	Definición y aspectos clínicos . . . . .	6
3.1.2	Biomarcadores y atrofia típica en MRI . . . . .	7
3.1.3	Escalas clínicas . . . . .	7
3.2	Aprendizaje automático y aprendizaje profundo en imagen médica . . . . .	7
3.2.1	Conceptos básicos de aprendizaje automático . . . . .	7
3.2.2	Redes neuronales y aprendizaje profundo . . . . .	8
3.2.3	Modelos generativos: GAN y WGAN . . . . .	10
3.2.4	Métricas de evaluación . . . . .	11
3.3	Estado del arte . . . . .	12
<b>4</b>	<b>Datos y preprocesamiento</b>	<b>14</b>
4.1	Descripción del dataset . . . . .	14
4.2	Pipeline de preprocesado . . . . .	16
4.2.1	Selección del tipo de MRI para cada sujeto . . . . .	16
4.2.2	Obtención del dataset de sujetos . . . . .	17
4.2.3	Skull Stripping . . . . .	18
4.2.4	Reorientación a RAS . . . . .	20
4.2.5	Redimensionado de los volúmenes . . . . .	21
4.2.6	Alineamiento y homogeneización del espacio . . . . .	21
4.3	División de conjuntos . . . . .	21

<b>5</b>	<b>Metodología de clasificación</b>	<b>23</b>
5.1	Arquitectura del modelo . . . . .	23
5.2	Entrenamiento del clasificador . . . . .	26
5.3	Métricas de evaluación . . . . .	28
<b>6</b>	<b>Metodología de generación</b>	<b>30</b>
6.1	Arquitectura del WGAN . . . . .	30
6.1.1	WGAN . . . . .	31
6.1.2	Arquitectura Final: WGAN crítico . . . . .	32
6.2	Estrategia de entrenamiento . . . . .	34
6.2.1	Funciones de pérdida . . . . .	34
6.2.2	Peso de reconstrucción y número de pasos del discriminador . . . . .	38
6.2.3	Optimizadores y scheduler de la tasa de aprendizaje . . . . .	39
6.2.4	Técnicas de estabilización adicionales . . . . .	40
6.3	Preprocesado y caché de volúmenes . . . . .	40
6.4	Uso previsto de las imágenes sintéticas . . . . .	41
6.4.1	Aumento de datos para el clasificador . . . . .	41
6.4.2	Análisis cualitativo y exploración visual . . . . .	41
6.4.3	Evaluación cuantitativa e impacto en el clasificador . . . . .	42
<b>7</b>	<b>Resultados</b>	<b>43</b>
7.1	Resultados del modelo de clasificación . . . . .	43
7.2	Clasificación binaria CN vs. Enfermedad . . . . .	48
7.2.1	Resultados con ResNet18 binaria . . . . .	49
7.3	Resultados del modelo de generación . . . . .	50
7.3.1	Modelo final de generación . . . . .	51
7.3.2	Diversidad morfológica de las muestras generadas . . . . .	53
7.3.3	Imágenes sintéticas finales . . . . .	55
7.3.4	Denoising de los cerebros sintéticos . . . . .	56
<b>8</b>	<b>Explicabilidad</b>	<b>59</b>
8.1	Métodos empleados . . . . .	59
8.2	Resultados de explicabilidad . . . . .	60
8.2.1	Explicabilidad en modelos con y sin datos sintéticos . . . . .	61
8.2.2	Comparación de los mapas Grad-CAM++ en modelos multiclase y binarios . . . . .	64
8.3	Implicaciones clínicas . . . . .	65
<b>9</b>	<b>Discusión</b>	<b>66</b>
9.1	Interpretación de resultados . . . . .	66

9.1.1	Comparación de modelos con y sin datos sintéticos . . . . .	68
9.2	Comparación con el estado del arte . . . . .	72
<b>10</b>	<b>Conclusiones y trabajo futuro</b>	<b>74</b>
10.1	Conclusiones . . . . .	74
10.2	Líneas de trabajo futuro . . . . .	76
	<b>Bibliografía</b>	<b>78</b>

# Capítulo 1

## Introducción y motivación

### 1.1 Contexto general del Alzheimer y diagnóstico por imagen

La enfermedad de Alzheimer es la causa más frecuente de demencia en personas de edad avanzada y se caracteriza por un deterioro progresivo de la memoria, las funciones ejecutivas y otras capacidades cognitivas, que acaba afectando de manera importante a la autonomía del paciente [1]. A nivel biológico, se asocia a la acumulación de placas de beta-amiloide, ovillos neurofibrilares de proteína tau y a un proceso de neurodegeneración que conlleva atrofia cerebral progresiva. Más allá de la descripción clínica, este deterioro se traduce en cambios muy concretos en la vida diaria: dificultades para reconocer lugares conocidos, problemas para organizar tareas sencillas o mantener una conversación, y una pérdida progresiva de iniciativa y de contacto con el entorno, que afecta no sólo a la persona enferma sino también a su familia y cuidadores [2].

En la práctica clínica, el diagnóstico se basa en la combinación de la historia del paciente, pruebas de memoria y otras funciones cognitivas, junto con diferentes pruebas complementarias [3]. Entre las pruebas cognitivas más utilizadas se encuentran escalas globales como el *Mini-Mental State Examination* (MMSE) y la *Clinical Dementia Rating* (CDR), que permiten cuantificar de forma sencilla el grado de deterioro y seguir su evolución en el tiempo [4], [5]. Entre estas herramientas, la neuroimagen tiene un papel muy importante, la resonancia magnética estructural (MRI) permite ver la forma y el tamaño de las distintas zonas del cerebro y detectar la pérdida de volumen en áreas que suelen verse afectadas en el Alzheimer, como el hipocampo y parte del lóbulo temporal [6].

La MRI se utiliza tanto para descartar otras posibles causas del deterioro cognitivo (por ejemplo, tumores o secuelas de ictus) como para apoyar el diagnóstico de Alzheimer

y seguir cómo evoluciona la enfermedad con el tiempo. Además, a partir de estas imágenes se pueden obtener medidas sencillas sobre el tamaño de ciertas estructuras cerebrales, que ayudan a los profesionales a tener una idea más objetiva del grado de afectación.

## 1.2 Motivación

En el ámbito de la imagen médica, y en particular en el estudio de la enfermedad de Alzheimer mediante resonancia magnética, es habitual trabajar con conjuntos de datos de tamaño limitado y con un marcado desequilibrio entre clases[7]. En muchos casos, las categorías más relevantes desde el punto de vista clínico, como determinados estadios tempranos de la enfermedad, están infrarrepresentadas, lo que dificulta que los modelos de aprendizaje profundo aprendan patrones fiables y tiende a sesgar las predicciones hacia las clases mayoritarias.

Esta escasez de datos y el desbalanceo de clases suponen un problema importante a la hora de desarrollar modelos de diagnóstico que sean robustos y generalizables. Para intentar mitigar estas limitaciones en la medida de lo posible, en los últimos años han aparecido modelos generativos capaces de sintetizar imágenes médicas realistas, lo que abre la puerta a utilizar estas imágenes sintéticas como técnica de aumento de dataset y equilibrio entre las clases. En esta línea, el presente trabajo se centra en entrenar un modelo de clasificación basado en aprendizaje profundo a partir de un conjunto finito de MRI estructurales y, posteriormente, en desarrollar un modelo generativo condicionado por clase que permita sintetizar nuevas imágenes; el objetivo final es evaluar si la incorporación de estas imágenes sintéticas al conjunto de entrenamiento contribuye a mejorar el rendimiento del clasificador y a reducir, al menos parcialmente, el impacto de la escasez y el desbalanceo de datos.

## 1.3 Planteamiento del problema y preguntas de investigación

El desarrollo de modelos de aprendizaje profundo para el diagnóstico de la enfermedad de Alzheimer a partir de MRIs tiene varias limitaciones prácticas: los conjuntos de datos suelen ser reducidos y desbalanceados, las fases tempranas de la enfermedad están poco representadas y muchos modelos funcionan como cajas negras difíciles de interpretar para el clínico.

Por ello, para resolver este problema general, este estudio se articula en torno a las

siguiente preguntas de investigación:

- ¿Es posible entrenar un clasificador 3D basado en aprendizaje profundo capaz de distinguir de forma fiable entre los diferentes estados cognitivos presentes en el conjunto de datos de MRI utilizado?
- ¿Puede un modelo generativo del tipo Wasserstein Generative Adversarial Network (WGAN) producir imágenes cerebrales sintéticas suficientemente realistas como para emplearlas como aumento de datos, especialmente en las clases minoritarias?
- ¿La incorporación de estas imágenes sintéticas al conjunto de entrenamiento mejora el rendimiento del modelo de clasificación frente al uso exclusivo de datos reales?
- ¿Qué regiones o estructuras cerebrales resultan más relevantes para las decisiones del clasificador y hasta qué punto son coherentes con el conocimiento clínico actual sobre la enfermedad de Alzheimer?

# Capítulo 2

## Objetivos y contribuciones

### 2.1 Objetivo general

El objetivo general de este estudio es desarrollar y evaluar un sistema basado en aprendizaje profundo para el análisis de resonancias magnéticas cerebrales tridimensionales (MRI 3D), cuyo núcleo sea un modelo de clasificación capaz de distinguir entre diferentes estados cognitivos. Dado que el número de datos disponibles es limitado y las clases están desbalanceadas, se entrena además un modelo generativo para sintetizar imágenes y añadirlas al conjunto de entrenamiento para analizar si mejora o no el rendimiento del clasificador.

### 2.2 Objetivos específicos

Para desarrollar este estudio se plantean tres objetivos principales. En primer lugar, diseñar, implementar y entrenar un clasificador tridimensional basado en redes neuronales convolucionales (CNN 3D) [8] capaz de distinguir entre diferentes estados cognitivos a partir de volúmenes de MRI preprocesados. En segundo lugar, implementar un modelo generativo de tipo WGAN que permita sintetizar volúmenes cerebrales realistas condicionados por clase, poniendo especial atención en las clases minoritarias del conjunto de datos. Por último, evaluar si la incorporación de imágenes generadas por el WGAN como técnica de aumento de datos contribuye a mejorar el rendimiento del clasificador frente al entrenamiento exclusivamente con datos reales. En tercer lugar, evaluar si la incorporación de imágenes sintéticas generadas por el WGAN como técnica de aumento de datos contribuye a mejorar el rendimiento del clasificador frente al entrenamiento exclusivamente con datos reales, e integrar además una capa de explicabilidad que permita interpretar tanto las decisiones del clasificador como las diferencias morfológicas entre clases.

## 2.3 Contribuciones

Las principales contribuciones de este trabajo se pueden resumir en varios aspectos complementarios. En primer lugar, se ha desarrollado un pipeline completo de preprocesado y clasificación 3D de resonancias magnéticas cerebrales aplicado a la enfermedad de Alzheimer, que incluye la definición de la arquitectura del modelo, el esquema de entrenamiento y la evaluación mediante métricas estándar en imagen médica. En segundo lugar, se ha implementado y entrenado un modelo WGAN capaz de generar imágenes cerebrales sintéticas condicionadas por clase, con el objetivo de abordar la escasez de datos y el desequilibrio de clases característicos de los conjuntos de datos de neuroimagen.

Además, se lleva a cabo un estudio experimental sistemático sobre el efecto del uso de datos sintéticos en el rendimiento del clasificador, comparando diferentes configuraciones de entrenamiento con y sin imágenes generadas para cuantificar la contribución real de este aumento de datos. Finalmente, se aplican técnicas de explicabilidad sobre el clasificador 3D para obtener mapas de relevancia que permiten relacionar las decisiones del modelo con regiones cerebrales descritas en la literatura sobre Alzheimer, discutiendo cómo estos resultados podrían resultar de utilidad en un contexto clínico.

# Capítulo 3

## Fundamentos teóricos

### 3.1 Enfermedad del Alzheimer

#### 3.1.1 Definición y aspectos clínicos

La enfermedad de Alzheimer es la causa más frecuente de demencia en personas mayores y se caracteriza por un deterioro progresivo de la memoria, el pensamiento y otras funciones cognitivas, que termina afectando a la autonomía y a la vida diaria del paciente. Los síntomas iniciales suelen incluir fallos de memoria recientes, dificultad para encontrar palabras, desorientación espacial leve y cambios de humor o de personalidad, que se hacen más evidentes a medida que la enfermedad avanza. En fases más avanzadas aparecen problemas importantes para comunicarse, reconocer a familiares, manejar tareas básicas y, finalmente, una dependencia casi completa de cuidadores.

Actualmente, la enfermedad del Alzheimer representa aproximadamente el 60–70 % de todos los casos de demencia a nivel mundial [9]. Se estima que el número de personas que viven con demencia se duplicará o incluso triplicará hacia el año 2050, en gran parte debido al envejecimiento de la población, lo que supone un incremento exponencial de la carga asistencial y económica asociada [10]. Este crecimiento no sólo tiene implicaciones sanitarias y de investigación, sino también un impacto humano profundo: detrás de cada diagnóstico hay una persona que va perdiendo progresivamente recuerdos, capacidades y autonomía, y un entorno familiar que debe adaptarse a cambios continuos en la conducta, la comunicación y el grado de dependencia del paciente.

### 3.1.2 Biomarcadores y atrofia típica en MRI

A nivel biológico, la enfermedad del Alzheimer se asocia a la acumulación anómala de placas de beta-amiloide y ovillos neurofibrilares de proteína tau, así como a un proceso de neurodegeneración que provoca pérdida progresiva de neuronas y sinapsis [11], [12]. Entre los biomarcadores de neuroimagen, uno de los más estudiados es la atrofia del hipocampo y de la corteza temporal medial, que suele aparecer de forma temprana y puede cuantificarse mediante resonancia magnética estructural. Estas medidas de volumen y grosor cortical se utilizan como apoyo al diagnóstico y para seguir la evolución de la enfermedad, aunque la atrofia hipocampal no es exclusiva del Alzheimer y debe interpretarse junto con otros datos clínicos y biomarcadores.

### 3.1.3 Escalas clínicas

El diagnóstico y la valoración de la severidad del Alzheimer se apoyan también en escalas clínicas estandarizadas. Una de las más utilizadas es el Mini-Mental State Examination (MMSE), que proporciona una puntuación global del rendimiento cognitivo y permite clasificar de forma aproximada la demencia en leve, moderada o grave [13]. Otra herramienta frecuente es la Clinical Dementia Rating (CDR), que valora diferentes dominios (memoria, orientación, juicio, actividades de la vida diaria, etc.) y ofrece una puntuación que refleja el estadio clínico del paciente [14]. La combinación de estas escalas con los hallazgos de neuroimagen y otros biomarcadores permite obtener una visión más completa del estado del paciente.

## 3.2 Aprendizaje automático y aprendizaje profundo en imagen médica

### 3.2.1 Conceptos básicos de aprendizaje automático

El aprendizaje automático (*machine learning*) es una rama de la inteligencia artificial en la que los modelos aprenden patrones a partir de datos, en lugar de estar programados explícitamente para cada tarea. En el caso del aprendizaje supervisado, que es el enfoque utilizado en este trabajo, el modelo recibe ejemplos de entrada junto con su etiqueta correcta (por ejemplo, una imagen de MRI y el diagnóstico y sujeto asociado) y aprende a predecir la etiqueta adecuada para nuevas entradas que no ha visto antes.

Para entrenar y evaluar estos modelos se suele dividir el conjunto de datos en tres

partes: entrenamiento, validación y test. El conjunto de entrenamiento se utiliza para ajustar los parámetros de la red; el de validación sirve para seleccionar hiperparámetros (como la tasa de aprendizaje, la arquitectura o el número de épocas) y detectar sobreajuste sin “mirar” directamente al test, y el conjunto de test se reserva para evaluar de forma final cómo se comporta el modelo sobre datos completamente nuevos. De este modo se obtiene una estimación más realista de la capacidad de generalización del modelo y se reduce el riesgo de que las métricas reflejen simplemente cómo de bien ha memorizado los ejemplos disponibles.

En el caso de redes neuronales profundas, el entrenamiento consiste en presentar sucesivamente lotes de imágenes al modelo, calcular la salida correspondiente y comparar esa salida con la etiqueta verdadera mediante una función de pérdida adecuada (por ejemplo, entropía cruzada para clasificación). Esta pérdida se deriva respecto a los pesos de la red mediante el algoritmo de retropropagación del error, y los pesos se actualizan utilizando un optimizador (como Stochastic Gradient Descent (SGD) o Adam) que aplica pequeños ajustes en la dirección que reduce dicha pérdida. Entre cada capa se utilizan funciones de activación no lineales (como Rectified Linear Unit (ReLU) o Leaky Rectified Linear Unit (LeakyReLU)), que permiten que la red aprenda relaciones complejas entre entrada y salida y evitan que el modelo se limite a aprender transformaciones lineales simples.

En el ámbito de la biomedicina y la imagen médica, estos métodos se aplican, entre otras cosas, a la detección automática de patologías, la clasificación de pacientes en diferentes estados de una enfermedad o la segmentación de estructuras anatómicas en distintas modalidades de imagen (MRI, Tomografía Computarizada (CT), rayos X, etc.). En el caso concreto de este trabajo, el aprendizaje supervisado se emplea para entrenar modelos capaces de distinguir entre distintos estados cognitivos a partir de resonancias magnéticas cerebrales.

### 3.2.2 Redes neuronales y aprendizaje profundo

Dentro del aprendizaje automático, el aprendizaje profundo (*deep learning*) es una rama que utiliza modelos llamados redes neuronales. Estos modelos están formados por muchas capas encadenadas, y cada capa transforma los datos de entrada, de manera que la red aprende por sí sola qué patrones son útiles para resolver una tarea, a partir de ejemplos etiquetados. Como se muestra en la Figura I, las neuronas artificiales están conectadas entre sí y cada conexión tiene un valor numérico asociado, denominado peso. El valor de salida de cada neurona se obtiene combinando (por ejemplo, mediante una suma ponderada) los valores de entrada con sus pesos correspondientes, lo que permite a la red ajustar esas conexiones durante el entrenamiento para mejorar sus predicciones.

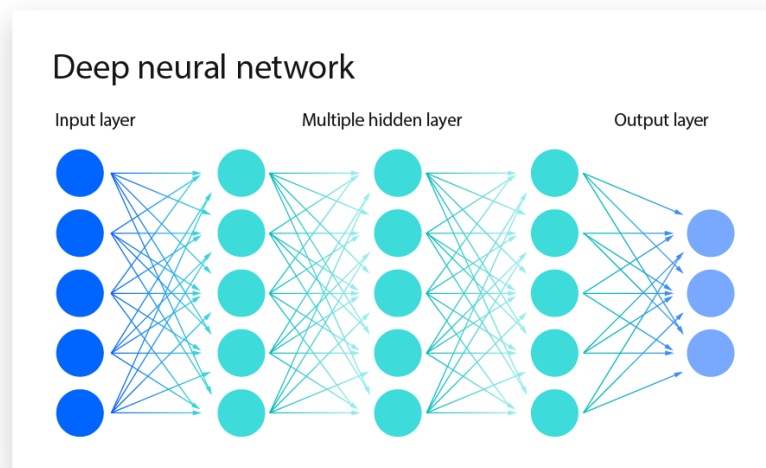


Figura 3.1: Estructura básica de una red neuronal[15]

Una red neuronal puede verse como un conjunto de neuronas artificiales conectadas entre sí, que reciben valores de entrada, los combinan y producen una salida. Durante el entrenamiento, la red ajusta automáticamente la intensidad de esas conexiones (pesos) para que sus predicciones se acerquen lo máximo posible a las etiquetas correctas, mediante un proceso de optimización de una función de coste que mide el error global del sistema. Este proceso se basa en el algoritmo de retropropagación del error, que permite actualizar los pesos capa a capa, y en el uso de funciones de activación no lineales en las neuronas (como ReLU o sigmoide), que hacen posible que la red aprenda patrones complejos a partir de los datos. Cuando se apilan muchas capas ocultas, hablamos de redes “profundas”, capaces de aprender representaciones cada vez más abstractas y ricas a partir de datos como imágenes, audio o texto.

En el caso de las imágenes, lo que se introduce en la red son los valores numéricos de sus píxeles (la intensidad de cada punto de la imagen). De forma simplificada, la primera capa de la red tiene tantas entradas como píxeles tenga la imagen: esta se “aplana” y se representa como un único vector de números, y a partir de ahí las siguientes capas van transformando esa información para extraer patrones útiles. En las resonancias magnéticas, que son imágenes tridimensionales, la idea es la misma, pero en lugar de tener solo alto y ancho se añade también la profundidad del volumen, de modo que el modelo trabaja con bloques 3D de voxels en lugar de simples píxeles 2D. Ocurre lo mismo que en las imágenes a color (RGB), cada píxel se representa mediante tres canales (rojo, verde y azul), de forma que la red procesa simultáneamente estos tres mapas de intensidad superpuestos en lugar de una única imagen en escala de grises [16].

Para trabajar de manera más eficiente con este tipo de datos se utilizan con frecuencia redes neuronales convolucionales (CNN), que aplican pequeños filtros que se desplazan por

la imagen y detectan patrones sencillos, como bordes o texturas, que después se combinan en estructuras más complejas [17]. En el caso de las imágenes médicas volumétricas, como las resonancias magnéticas cerebrales, es habitual emplear variantes tridimensionales de estas redes (CNN 3D), que aplican los filtros en las tres dimensiones y trabajan directamente con el volumen completo en lugar de analizar cada corte por separado [18]. A medida que se encadenan capas de convolución y de *pooling*, la red va reduciendo progresivamente la dimensionalidad de los datos pero conservando la estructura espacial de los píxeles o vóxeles más relevante, hasta llegar a una o varias capas totalmente conectadas que operan sobre una representación compacta del volumen y permiten realizar finalmente la clasificación.

### 3.2.3 Modelos generativos: GAN y WGAN

Los modelos generativos intentan aprender cómo se distribuyen los datos de entrada para poder crear nuevos ejemplos que se parezcan a los reales. Dentro de esta familia, las redes generativas antagónicas (GAN)[19] utilizan dos redes que se entrenan a la vez y compiten entre sí: un generador, que produce imágenes sintéticas a partir de ruido aleatorio, y un discriminador, que intenta decidir si cada imagen que ve es real (procedente del conjunto de datos) o falsa (producida por el generador). Durante el entrenamiento, el generador se actualiza para engañar al discriminador, mientras que el discriminador se entrena para distinguir cada vez mejor entre ejemplos reales y generados, de manera que ambos modelos mejoran conjuntamente.

Este proceso se formula como un juego *minimax*: el discriminador intenta maximizar una función de coste que recompensa clasificar correctamente las imágenes reales y penaliza confundir las falsas, mientras que el generador intenta minimizar esa misma función haciendo que las imágenes sintéticas resulten indistinguibles de las reales[19], [20]. En la formulación original de Goodfellow, este juego se expresa como

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{datos}}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))],$$

donde  $D(x)$  es la probabilidad de que la muestra  $x$  sea real y  $G(z)$  es la imagen generada a partir de un vector de ruido  $z$ [19]. En la práctica, este esquema puede presentar problemas de inestabilidad, como gradientes que desaparecen cuando el discriminador se vuelve demasiado fuerte o el fenómeno de *mode collapse*, en el que el generador solo produce un conjunto limitado de ejemplos distintos[21]. Debido a estas limitaciones se han propuesto muchas variantes; entre ellas, las Wasserstein GAN (WGAN), que son las que se emplean en este trabajo, sustituyen la función de pérdida clásica basada en la divergencia de Jensen–Shannon por una aproximación de la distancia de Wasserstein entre la distribución

real y la generada[22]. Esta distancia tiene la ventaja de seguir proporcionando gradientes informativos incluso cuando las distribuciones apenas se solapan, evitando que el generador deje de aprender cuando el discriminador es demasiado preciso. Versiones mejoradas como WGAN-GP incorporan una penalización de gradiente (*gradient penalty*) en lugar del recorte de pesos clásico, lo que impone de forma más estable la condición de Lipschitz al discriminador y mejora la calidad de las imágenes generadas. En este contexto, el discriminador se denomina *critic*, ya que en lugar de producir una probabilidad de que la imagen sea real devuelve una puntuación escalar relacionada con esa distancia, y es habitual entrenarlo varias veces por cada actualización del generador (parámetro  $n_{critic}$ ), con el objetivo de que se aproxime mejor a su solución óptima y proporcione gradientes más fiables al generador[23], [24].

### 3.2.4 Métricas de evaluación

Para evaluar modelos de clasificación en imagen médica se emplean distintas métricas derivadas de la matriz de confusión. Estas métricas han sido vitales durante el entrenamiento de los modelos, ya que gracias a ellas se ha podido ajustar los hiperparámetros correctamente para poder obtener el mejor modelo posible. Existen muchas, pero para la evaluación de este proyecto se han utilizado principalmente la exactitud, la sensibilidad, la especificidad, la precisión, la F1 y el área bajo la curva ROC (AUC). El *exactitud* mide el porcentaje total de aciertos, pero puede resultar engañosa cuando las clases están desbalanceadas, por lo que suele complementarse con medidas como la sensibilidad (recall de la clase positiva), la especificidad y la precisión. La F1 combina precisión y sensibilidad en un único valor, útil cuando interesa equilibrar ambas, mientras que la curva ROC y el área bajo la curva (AUC) permiten valorar el rendimiento del modelo a distintos umbrales de decisión y son muy utilizadas en problemas biomédicos.

En el caso de los modelos generativos de imágenes, como las GAN y WGAN utilizadas en este trabajo, se emplean métricas adicionales orientadas a valorar tanto la calidad visual como la diversidad de las imágenes sintéticas[25], [26]. Entre las más habituales se encuentran el *Inception Score* (IS), que mide si las imágenes generadas son claras y variadas analizando las predicciones de una red de clasificación preentrenada, y la *Fréchet Inception Distance* (FID), que compara estadísticamente las representaciones de alto nivel de imágenes reales y generadas: valores de FID más bajos indican que las distribuciones son más similares[27]. En imagen médica es frecuente complementar estas métricas con indicadores de calidad de imagen como la relación señal-ruido (SNR), el *peak signal-to-noise ratio* (PSNR) y el índice de similitud estructural (SSIM), que permiten cuantificar hasta qué punto se preservan la intensidad y las estructuras anatómicas de interés entre la imagen original y la generada o denoised[28].

### 3.3 Estado del arte

En los últimos años se han propuesto numerosos modelos de aprendizaje profundo para la clasificación de la enfermedad de Alzheimer a partir de resonancias magnéticas estructurales. La mayoría de trabajos emplean redes convolucionales 2D o 3D entrenadas sobre bases de datos públicas como ADNI para distinguir entre sujetos con Alzheimer y controles sanos, alcanzando con frecuencia *accuracies* superiores al 80 % y valores de AUC por encima de 0,9 en tareas binarias AD (Alzheimer’s Disease) vs. CN (Cognitively Normal)[29]. Cuando se aborda la clasificación multiclase de distintos estadios (CN, MCI (Mild Cognitive Impairment), AD) o la predicción de la conversión de MCI a AD, el rendimiento suele disminuir y aparecen con más claridad problemas de sobreajuste, especialmente en estudios con tamaños muestrales reducidos, protocolos de validación poco estrictos o ausencia de validación externa. A pesar de que se han explorado arquitecturas cada vez más complejas (por ejemplo, modelos 3D, redes de atención o combinaciones con transformadores) y el uso de información multimodal, la variabilidad entre estudios en cuanto a preprocesado, división de datos y métricas dificulta la comparación directa de resultados[30].

En paralelo, el uso de modelos generativos, especialmente difusión models, GAN y sus variantes, ha ganado relevancia en imagen médica como herramienta para sintetizar datos adicionales y paliar la escasez y el desbalance de los conjuntos de entrenamiento. Diversos trabajos han demostrado que es posible generar imágenes realistas de MRI, CT o rayos X y utilizarlas como *aumento de datos* para mejorar el rendimiento de modelos de clasificación o segmentación, incluyendo aplicaciones concretas en el diagnóstico de Alzheimer. En este contexto se emplean tanto métricas clásicas de calidad de imagen (PSNR, SSIM, SNR) como medidas específicas para modelos generativos, tales como el *Inception Score* (IS) y la *Fréchet Inception Distance* (FID), que permiten cuantificar la similitud entre la distribución de las imágenes sintéticas y la de las imágenes reales [31]. Además, el uso de datos sintéticos contribuye a reducir algunos problemas de privacidad y de intercambio de datos entre centros, ya que estas imágenes no pertenecen a ningún paciente real.

Sin embargo, los enfoques existentes presentan todavía limitaciones importantes. Por un lado, muchos estudios se basan en muestras pequeñas o en un único centro, lo que favorece el sobreajuste y dificulta la generalización a otros hospitales, escáneres o poblaciones [29], [32]. Por otro, la mayoría de modelos funcionan como “cajas negras” con un grado limitado de explicabilidad, lo que complica su aceptación en entornos clínicos donde es importante entender qué regiones o patrones de la imagen sustentan la decisión del modelo [33]. En el ámbito generativo, aunque las GAN y WGAN producen imáge-

nes visualmente plausibles, métricas como FID o IS no siempre garantizan que se hayan preservado las estructuras anatómicas relevantes ni que las muestras generadas aporten realmente información nueva (riesgo de memorizar el conjunto de entrenamiento o de sufrir un colapso del modelo) [34], [35]. Todo ello pone de manifiesto la necesidad de desarrollar modelos que, además de obtener buenos resultados en las métricas habituales, se evalúen con protocolos rigurosos, tengan en cuenta la variabilidad entre centros y aporten un mayor grado de interpretabilidad y control sobre la calidad de las imágenes sintéticas.

# Capítulo 4

## Datos y preprocesamiento

### 4.1 Descripción del dataset

Este trabajo utiliza el conjunto de datos *Alzheimer’s Disease Neuroimaging Initiative* (ADNI)[36], un estudio multicéntrico que recoge imágenes de resonancia magnética (MRI), información clínica y diversos biomarcadores de sujetos con diferente estado cognitivo. Para el desarrollo del modelo se selecciona una muestra de 882 sujetos, divididos en tres grupos: 458 controles cognitivamente normales (CN), 231 pacientes con deterioro cognitivo leve (MCI) y 193 pacientes con diagnóstico de enfermedad de Alzheimer (AD). En la Figura 4.1 se muestran ejemplos de imágenes T1 representativas de cada uno de estos grupos diagnósticos (CN, MCI y AD), donde se aprecian las diferencias de atrofia cerebral asociadas a la progresión de la enfermedad. En las imágenes se aprecia que la enfermedad del Alzheimer (Grupo AD) se puede intuir la diferencia debido a la atrofia que reciben regiones estructurales del cerebro, mientras que la distinción entre CN y MCI resulta mucho más sutil a simple vista; precisamente, uno de los objetivos de este trabajo es que el modelo de clasificación aprenda a capturar esas diferencias tempranas difíciles de detectar de forma visual directa.

La obtención y selección de las imágenes no es un proceso trivial. El repositorio ADNI incluye distintos tipos de estudios y secuencias de resonancia magnética, por lo que es necesario revisar con detalle la documentación disponible para identificar qué series eran comparables entre sí y garantizar que todas las imágenes utilizadas en el proyecto fuesen del mismo tipo. Por todo ello, se toma la decisión de trabajar con imágenes ponderadas en T1, ya que proporcionan una representación estructural clara de la anatomía cerebral y permiten visualizar con buena resolución las diferencias de volumen y forma entre regiones como el hipocampo y la corteza temporal, que son especialmente relevantes en la detección de la enfermedad de Alzheimer.

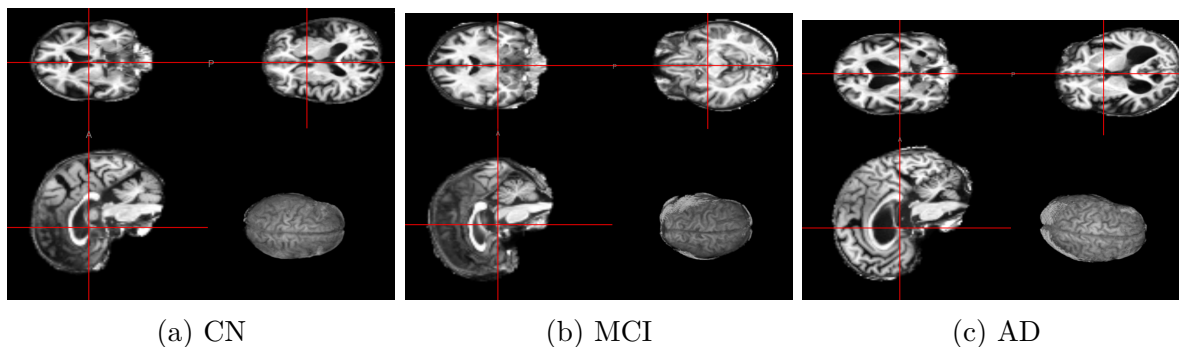


Figura 4.1: Ejemplos de imágenes T1 estructurales de sujetos a) cognitivamente normales (CN), b) con deterioro cognitivo leve (MCI) y c) con enfermedad de Alzheimer (AD) procedentes del dataset ADNI.

Otras secuencias, como las que ponen más énfasis en el flujo sanguíneo o en el contenido de líquido (por ejemplo, secuencias T2 o angiográficas), resultan muy útiles para estudiar vasos, edema o procesos vasculares, pero no son tan adecuadas cuando el objetivo principal es analizar cambios estructurales sutiles en el cerebro, como es nuestro caso. Por este motivo, y para mantener un conjunto de datos homogéneo y centrado en la detección de atrofia y otros patrones, en este trabajo se ha decidido limitar el estudio a las imágenes T1 estructurales.

Dado que las imágenes proceden de diferentes proveedores, escáneres o fabricantes (por ejemplo, Philips o Siemens), poseen características que deben ser homogeneizadas como variaciones en parámetros de adquisición y contraste de imagen. Esto ha hecho necesario aplicar un preprocesado cuidadoso y, en algunos casos, ligeramente diferente según el origen del estudio, con el objetivo de reducir estas diferencias técnicas y disponer de volúmenes lo más homogéneos posible antes de entrenar los modelos de aprendizaje profundo.

Además, en la plataforma de ADNI las imágenes de MRI pueden descargarse en distintos formatos según el nivel de procesado: *original*, *processed* y *pre-processed*. Las imágenes *original* corresponden a los datos tal y como salen del escáner, mientras que las *processed* han pasado por algunos pasos básicos de corrección y estandarización definidos por el propio estudio. En este trabajo se decide por utilizar las imágenes *pre-processed*, que ya incorporan correcciones adicionales específicas para T1 (como gradwarp e intensidad normalizada) y se distribuyen en formato volumétrico *.nii*, lo que facilita cargarlas directamente como matrices tridimensionales y reduce la carga de implementar desde cero todo el pipeline de correcciones.

## 4.2 Pipeline de preprocesado

El preprocesado de las resonancias magnéticas ha supuesto la parte más crítica y crucial de este proyecto, hasta el punto de condicionar directamente la calidad de los resultados obtenidos con los modelos de *aprendizaje profundo*. Antes de comenzar a hacer los modelos de clasificación y generación, ha sido necesario invertir una cantidad considerable de tiempo en entender cómo están organizados los datos de ADNI, qué tipos de imágenes están disponibles para cada sujeto y qué nivel de procesado necesitaba cada tipo de imagen.

### 4.2.1 Selección del tipo de MRI para cada sujeto

Para cada sujeto del estudio ADNI no se dispone de una única resonancia T1, sino de varias versiones de la misma adquisición con diferentes niveles de procesado aplicados por el *MRI Core*. En la base de datos aparecen descripciones de imagen como, por ejemplo, “*MPRAGE; original*”, “*MPR-R; GradWarp; B1 Correction; N3; Scaled*” o “*Spatially Normalized, Masked and N3 corrected T1 image*”, entre otras. Esto obligó a definir una jerarquía clara para elegir, en cada sujeto, qué tipo de imagen utilizar, ya que algunos de estos tipos de MRIs al estar ya procesados facilitan la detección de patrones para nuestro modelo de clasificación.

De manera general, las imágenes etiquetadas con **N3** son volúmenes a los que se ha aplicado una corrección de no uniformidad de intensidad (*N3 bias field correction*), que variaciones de brillo que no dependen del cerebro del paciente, sino del propio escáner. Cuando además incluyen la marca **Scaled**, significa que se ha aplicado una corrección de escala y geometría basada en un fantoma de calibración, destinada a compensar pequeñas distorsiones en el tamaño y la forma del volumen. Por su parte, las imágenes descritas como “**Spatially Normalized, Masked and N3 corrected T1**”, además de la corrección N3, vienen alineadas a un mismo espacio de referencia y con una máscara que deja solo el cerebro, eliminando la mayor parte del cráneo y otros tejidos externos. Como se puede ver en la figura 4.2 se ilustra cómo cambia la apariencia de la MRI de un mismo paciente según el tipo de procesado aplicado.

En este trabajo se plantea una jerarquía clara a la hora de elegir qué versión de la MRI utilizar en cada sujeto, con el objetivo de aprovechar al máximo el preprocesado que ya ofrece ADNI y reducir el trabajo adicional necesario.

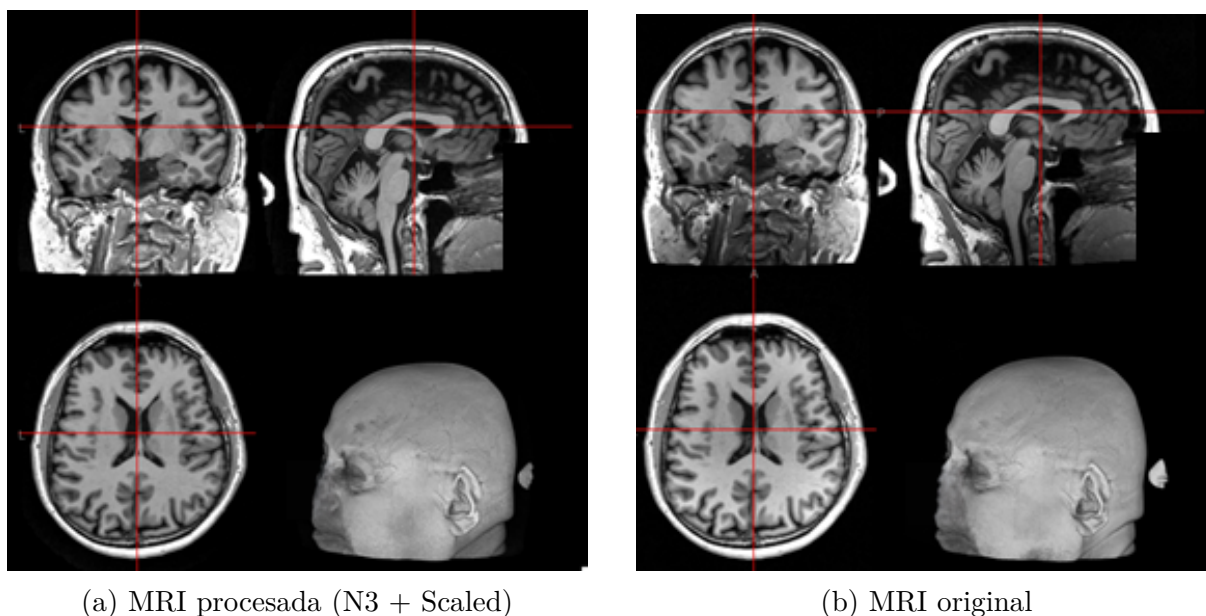


Figura 4.2: Comparación entre la imagen T1 original y la imagen T1 procesada con corrección N3 y *Scaled* para un mismo sujeto.

## 4.2.2 Obtención del dataset de sujetos

A partir de las tablas de ADNI que describen cada imagen individual (*Image Data ID*, modalidad, descripción del procesado, etc.), es necesario extraer y reorganizar la información asociada para construir el dataset de trabajo. En estos ficheros originales aparecen varias filas por sujeto, una por cada versión de la misma adquisición (por ejemplo, “MPR; ; N3; Scaled”, “MPR-R; ; N3; Scaled”, etc.), con columnas como *Subject*, *Group*, *Sex*, *Age* y *Acq Date*.

Mediante funciones de la librería `pandas` de Python se lleva a cabo un proceso de limpieza y filtrado para eliminar duplicidades por sujeto y visita, conservando una única fila representativa para cada participante incluido en el estudio. El resultado final es la tabla 4.1, más compacta con las columnas *Subject*, *Group*, *Sex*, *Age* y *Acq Date*, que se utiliza como base para enlazar las imágenes seleccionadas con sus correspondientes etiquetas clínicas.

Subject	Group	Sex	Age	Acq Date
136_S_1227	MCI	F	65	02/21/2007
136_S_0579	MCI	F	66	07/10/2006
136_S_0429	MCI	M	63	06/27/2006
136_S_0426	AD	M	80	05/30/2006
136_S_0196	CN	F	78	05/01/2006

Tabla 4.1: Ejemplo del dataset final de sujetos generado a partir de las tablas originales de imágenes de ADNI.

### 4.2.3 Skull Stripping

Tanto para la tarea de clasificación como para la generación de nuevos cerebros sintéticos, es necesario trabajar únicamente con el tejido cerebral, ya que es en el propio cerebro donde se manifiestan los cambios estructurales asociados a la enfermedad, mientras que el cráneo y otros tejidos externos no aportan información directa para detectar e identificar el diagnóstico. Como hemos mencionado anteriormente, debido a la heterogeneidad del conjunto de datos MRI, algunas imágenes T1 ya incluían una máscara binaria de cerebro generada en pasos de preprocesado previos, mientras que otras no incluían dicha máscara. Por ello, se llevaron a cabo dos preprocesamientos diferentes: en los casos en los que la máscara estaba disponible, el *skull stripping* se realizó de forma directa multiplicando voxel a voxel la imagen original por su máscara asociada, lo que permitió eliminar el cráneo y los tejidos no cerebrales y obtener un volumen de cerebro limpio sin modificar la intensidad de la región que nos interesa.

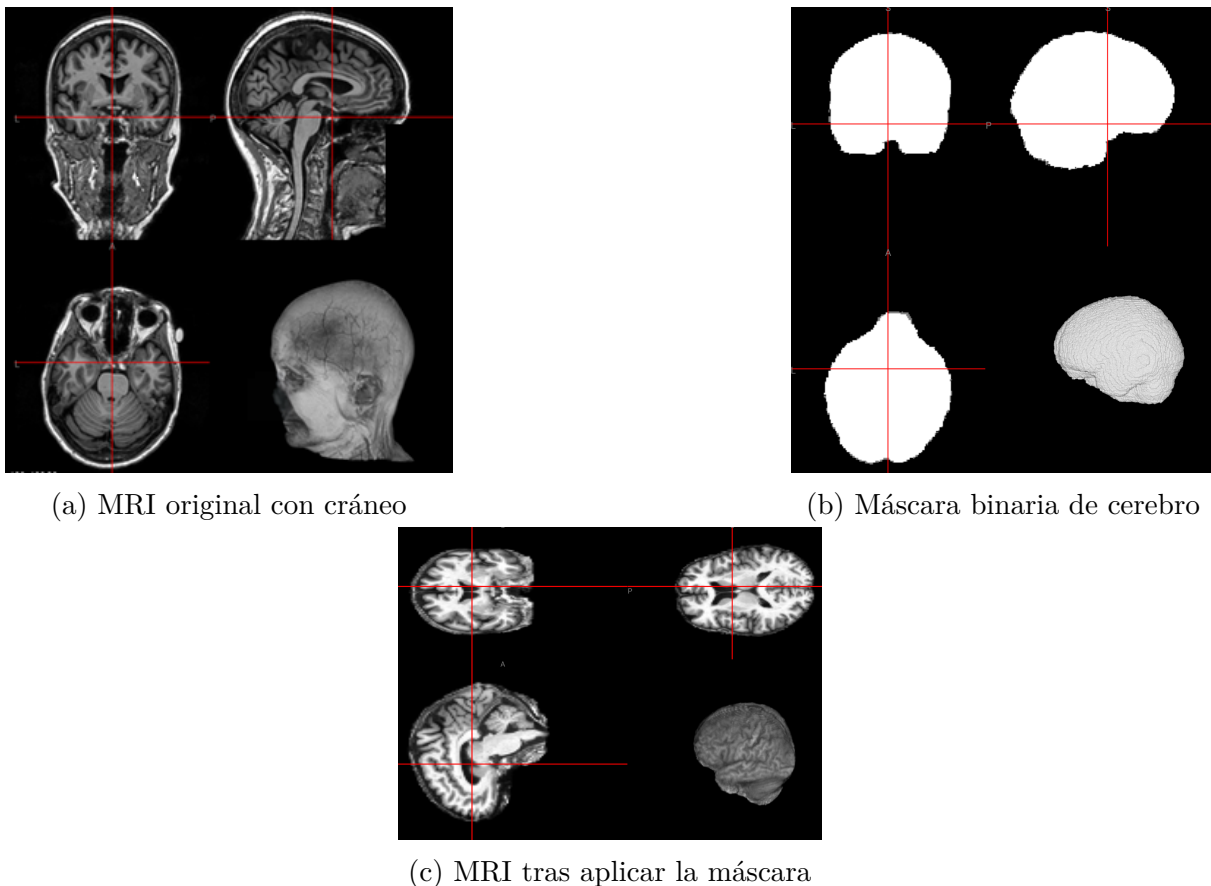


Figura 4.3: Ejemplo del proceso de *skull stripping* mediante operación voxel a voxel: a partir de a) la imagen T1 original y de b) su máscara binaria de cerebro se obtiene c) un volumen libre de cráneo y otros tejidos no cerebrales.

En la Figura 4.3 se ilustra este proceso paso a paso, desde la imagen T1 original hasta la obtención del volumen final libre de cráneo tras aplicar la máscara de cerebro, mostran-

do visualmente cómo se suprimen los tejidos no cerebrales y se conserva únicamente la anatomía cerebral relevante para el modelo. En aquellas imágenes T1 que no disponían de una máscara binaria asociada fue necesario generar dicha máscara de forma automática. Para ello se empleó *SynthStrip*[37], un modelo de *skull stripping* basado en aprendizaje profundo publicado junto con su implementación en GitHub[38], diseñado específicamente para extraer el cerebro de forma robusta a partir de imágenes de muy distinto contraste y protocolo de adquisición. En nuestro pipeline, *SynthStrip* se ejecutó sobre cada volumen T1 de entrada y devolvió dos salidas: por un lado, la imagen ya libre y con el cráneo y otros tejidos no cerebrales eliminados, y por otro, la máscara binaria de cerebro correspondiente, que se almacenó para mantener el mismo formato que en el caso de las imágenes que ya venían enmascaradas de origen.

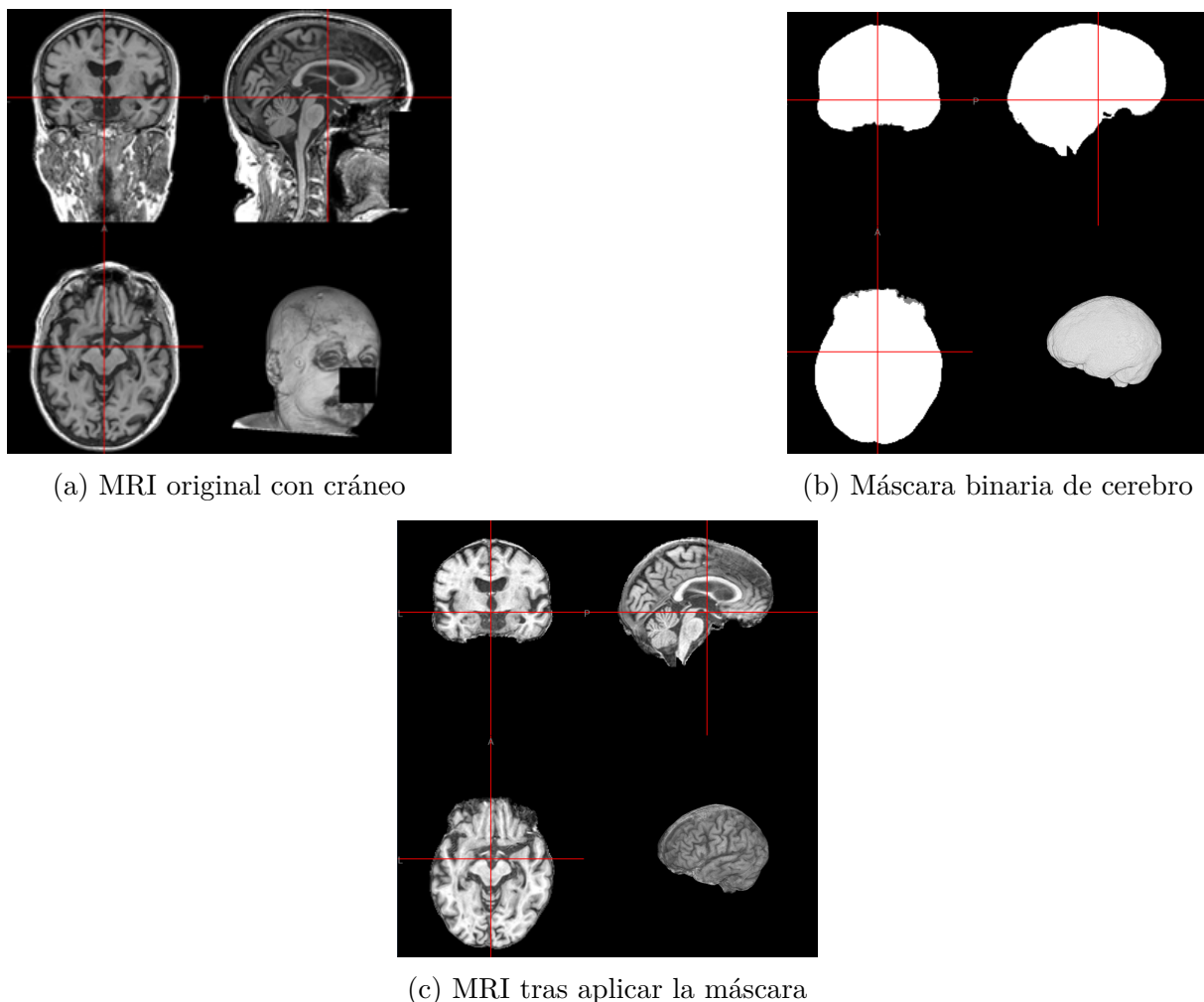


Figura 4.4: Ejemplo del proceso de *skull stripping* con el modelo synthtrip: a partir de la imagen T1 original se obtiene un volumen libre de cráneo y otros tejidos no cerebrales y su máscara correspondiente.

De este modo, tanto las imágenes con máscara preexistente como aquellas procesadas con *SynthStrip* terminaron representadas por un par (MRI sin cráneo, máscara de cerebro), garantizando un tratamiento homogéneo del *skull stripping* en todo el conjunto de da-

tos y evitando que las diferencias entre centros o protocolos afectaran al entrenamiento de los modelos posteriores. En la Figura 4.4 se ilustra gráficamente este procedimiento, mostrando la imagen T1 original, la máscara binaria estimada y el resultado final tras la extracción de cráneo, donde únicamente se conserva el tejido cerebral relevante para las tareas de clasificación y generación.

#### 4.2.4 Reorientación a RAS

Antes de entrenar tanto el modelo de clasificación como el modelo generativo resulta necesario asegurar que todas las imágenes compartan la misma orientación espacial. Para ello, todos los volúmenes T1 se reorientan al sistema estándar RAS (*Right–Anterior–Superior*), en el que los ejes de coordenadas aumentan hacia la derecha, hacia delante y hacia arriba del sujeto, respectivamente.

Esta homogeneización de la orientación evita que los modelos aprendan patrones contradictorios debidos únicamente a rotaciones o volteos de las imágenes, lo que en el caso del modelo generativo se traduciría en cerebros sintéticos con mitades mal alineadas o invertidas. Al trabajar con todos los sujetos en la misma convención RAS, la red solo se enfrenta a variaciones anatómicas relevantes y no a cambios arbitrarios de orientación que podrían confundir el proceso de aprendizaje. En la Figura 4.5 se muestra un ejemplo de dos MRIs del mismo sujeto con orientaciones diferentes, donde se aprecia cómo estos cambios pueden alterar de forma notable la apariencia global del volumen pese a corresponder al mismo cerebro.

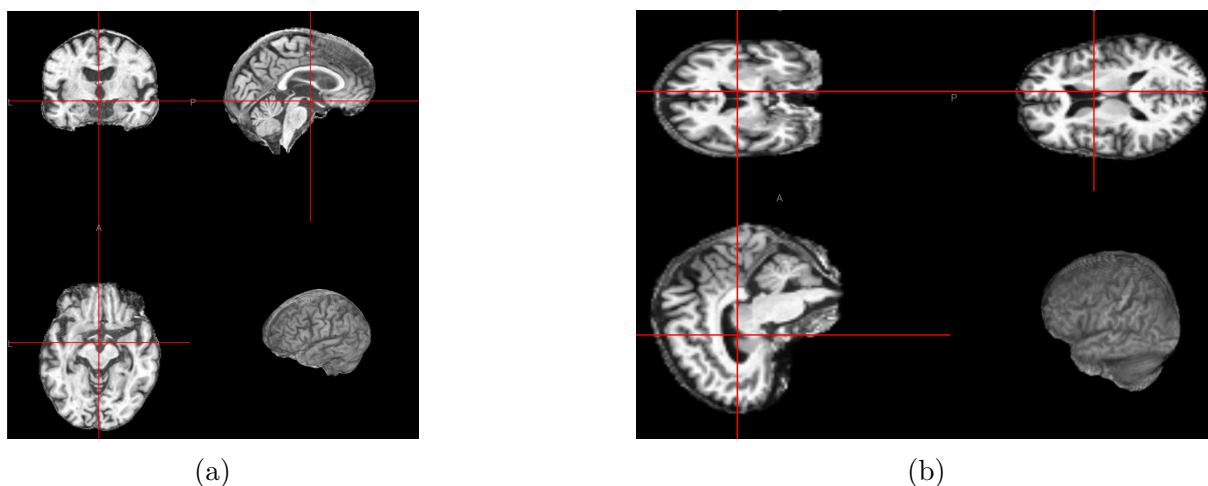


Figura 4.5: Diferencias entre dos MRIs con a) Orientación RAS y b) Orientación LAS

### 4.2.5 Redimensionado de los volúmenes

Una vez corregida la orientación, los volúmenes T1 seguían presentando tamaños de matriz y resoluciones diferentes según el centro de adquisición y la máquina con la que se realizó el escaner. Sin embargo, las redes neuronales 3D empleadas en este trabajo requieren que todas las entradas tengan exactamente las mismas dimensiones espaciales para poder formar lotes de entrenamiento y definir correctamente las capas convolucionales.

Por este motivo, cada volumen se re-muestrea a un tamaño fijo de  $128 \times 128 \times 128$  voxels, utilizando interpolación para ajustar la información original a la nueva rejilla uniforme. Este paso garantiza que todos los cerebros se representen en un espacio discreto común y permite que tanto el modelo de clasificación como el generativo reciban entradas compatibles y comparables entre sí.

### 4.2.6 Alineamiento y homogeneización del espacio

Además del tamaño del volumen, es importante que la posición relativa del cerebro dentro de la imagen sea lo más similar posible entre sujetos. Para ello, se aplica un alineamiento adicional que centra y encuadra los cerebros de manera consistente, de forma que quede aproximadamente en la misma región del campo de visión y las zonas de fondo presenten un aspecto homogéneo.

Este alineamiento reduce diferencias triviales de traslación y recorte entre estudios, evitando que los modelos aprendan a explotar artefactos de borde o variaciones del fondo en lugar de centrarse en la morfología cerebral. Como resultado, todas las imágenes empleadas en el entrenamiento comparten no solo la misma orientación y tamaño, sino también un encuadre coherente, lo que facilita la comparación directa entre sujetos y mejora la estabilidad del aprendizaje.

## 4.3 División de conjuntos

Para evaluar de forma fiable el rendimiento de los modelos se divide el conjunto de datos en tres particiones no solapadas: entrenamiento (*train*), validación (*validation*) y prueba (*test*). El conjunto de entrenamiento se utiliza para ajustar los parámetros de la red neuronal, el de validación para seleccionar hiperparámetros y decidir criterios de parada temprana, y el conjunto de prueba se mantiene completamente reservado hasta el final del experimento, empleándose únicamente para estimar el rendimiento final del modelo en sujetos no vistos.

Dado que las clases (CN, MCI, AD) están desbalanceadas, la división se realiza de manera estratificada por clase para preservar, en la medida de lo posible, la misma proporción de diagnósticos en las tres particiones. Además, la división se hace a nivel de sujeto y no a nivel de imagen, de forma que todas las resonancias de un mismo paciente queden siempre en el mismo conjunto (ya fuera entrenamiento, validación o prueba) y nunca repartidas entre varios. Esta estratificación por sujeto/clase evita que el modelo vea en entrenamiento imágenes prácticamente idénticas a las de prueba (por ejemplo, distintos cortes o visitas del mismo paciente), lo que podría inflar artificialmente las métricas debido a *data leakage*, es decir, a la presencia de información del conjunto de prueba “filtrada” implícitamente en el conjunto de entrenamiento [39], [40].

Para prevenir fugas de información adicionales, todas las transformaciones dependientes de los datos (normalización de intensidades, estandarización de variables tabulares, etc.) se ajustan exclusivamente sobre el conjunto de entrenamiento, y los parámetros resultantes se aplican después, sin volver a recalcularlos, a las particiones de validación y prueba. De esta forma se garantiza que ninguna información estadística procedente de los sujetos de validación o prueba influya en el proceso de entrenamiento o en la selección de hiperparámetros, manteniendo una separación estricta entre desarrollo del modelo y evaluación final.

Finalmente, en la Figura 4.6 se muestra de forma esquemática el flujo completo de preprocesamiento aplicado a las imágenes T1, desde la descarga de los datos en ADNI hasta la generación de los conjuntos de entrenamiento, validación y prueba utilizados por los modelos de aprendizaje profundo.

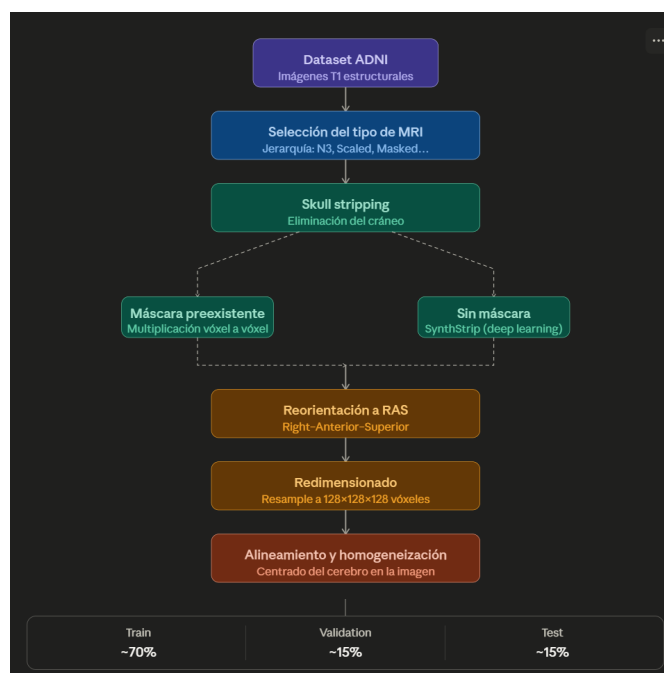


Figura 4.6: Flujo Preprocesamiento

# Capítulo 5

## Metodología de clasificación

### 5.1 Arquitectura del modelo

La definición de la arquitectura final requiere el análisis previo de diversas alternativas de complejidad incremental, lo que fundamenta la solidez técnica de la solución adoptada. En un primer momento se implementa una CNN 3D relativamente sencilla para comprobar que el *pipeline* de preprocesado y la división de datos funcionan correctamente. Esta red consta de varios bloques de convolución tridimensional con kernels  $3 \times 3 \times 3$ , seguidos de capas de *batch normalization*, activaciones ReLU y operaciones de *max-pooling*, terminando en una o varias capas totalmente conectadas para realizar la clasificación en tres clases (CN, MCI y AD). Aunque este modelo básico permite validar el flujo de datos, sus resultados en validación son modestos y ponen de manifiesto que se necesita una arquitectura más potente.

Como siguiente paso se prueba una arquitectura 3D específica para la enfermedad de Alzheimer, diseñada para trabajar con volúmenes T1 de ADNI y capturar patrones de atrofia característicos de la enfermedad [41]. Este modelo, más profundo que la CNN inicial, mejora la precisión hasta situarla en torno al 65 % en validación, lo que confirma que la información presente en las imágenes permite aprender un clasificador razonable pero que todavía queda margen de mejora, especialmente en la distinción entre sujetos CN y MCI.

### ResNet 3D: aprendizaje residual

La arquitectura final adoptada en este trabajo se basa en una ResNet 3D. Una ResNet (Residual Network) es un tipo de red neuronal profunda diseñada para facilitar el en-

trenamiento de modelos muy profundos mediante el uso de conexiones residuales, que permiten que la información fluya directamente entre capas no consecutivas [42]. Este enfoque ayuda a mitigar problemas como la degradación del rendimiento y el desvanecimiento del gradiente a medida que aumenta la profundidad de la red. Inspirada en la propuesta original de He et al. para clasificación de imágenes [42] y adaptada al contexto de volúmenes de resonancia magnética cerebral, esta arquitectura introduce mejoras clave frente a las CNN tradicionales. A diferencia de una CNN “clásica”, en la que cada bloque transforma la imagen de forma secuencial, sin reutilizar explícitamente la información de las capas anteriores, las redes residuales añaden conexiones de atajo (*skip connections*) que permiten que la información de entrada de un bloque se sume directamente a su salida, tal y como se ilustra en la Figura 5.1.

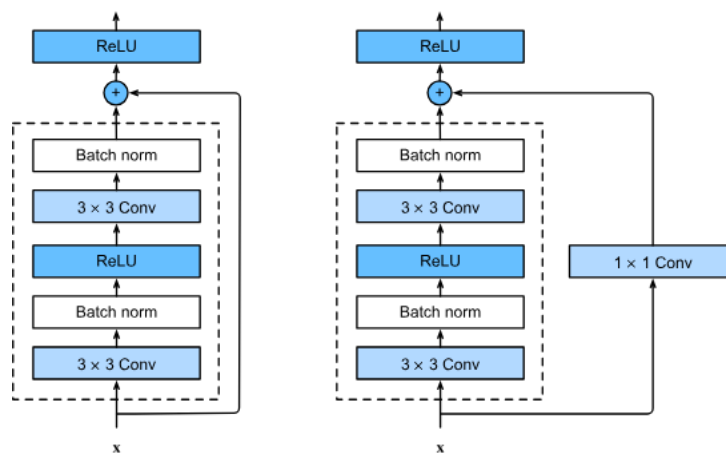


Figura 5.1: Representación esquemática de un bloque residual 3D con skip connection [43].

Intuitivamente, estas conexiones hacen que la red no vaya “perdiendo” la información original a medida que atraviesa capas. La salida de cada bloque no depende solo de las nuevas transformaciones que se aplican, sino que también incluye directamente la información que entra al bloque. De este modo, los patrones y características importantes de la imagen se mantienen presentes a lo largo de toda la red, incluso cuando esta es muy profunda. Si un bloque no necesita modificar nada relevante, la *skip connection* le permite dejar pasar casi intacta la información; y, sin embargo, si existen aspectos que mejorar, el bloque aprende un ajuste adicional más pequeño y específico. Esta idea hace que las redes profundas sean más estables de entrenar y ayuda a evitar problemas típicos cuando se añaden muchas capas, como que el entrenamiento se vuelva inestable o que las últimas capas acaben “deshaciendo” lo que aprenden las anteriores.

En nuestro caso se implementa una familia de ResNets tridimensionales parametrizadas por el número de bloques en cada etapa (**resnet18**, **resnet34**, **resnet50**, **resnet101**, etc.), utilizando tanto bloques residuales simples (*BasicBlock*) como bloques de tipo

*Bottleneck.* Los bloques *Bottleneck* introducen una estructura de tres convoluciones ( $1 \times 1 \times 1$ ,  $3 \times 3 \times 3$ ,  $1 \times 1 \times 1$ ) que primero reduce el número de canales, aplica la convolución y operación costosa en un espacio mucho más comprimido y finalmente vuelve a expandir la dimensionalidad, lo que permite construir redes mucho más profundas sin que el número de parámetros crezca de forma incontrolada.

La idea es similar a la que se muestra de forma esquemática en la Figura 5.2: la información de entrada pasa por una zona estrecha (*bottleneck*) donde se fuerza a la red a concentrar en pocas neuronas las características realmente importantes, y después se vuelve a expandir la representación. En el caso de la ResNet, ese papel de “cuello de botella” lo desempeñan las convoluciones  $1 \times 1 \times 1$  que comprimen y luego reexpanden los canales, haciendo más eficiente el cálculo sin perder la información relevante.

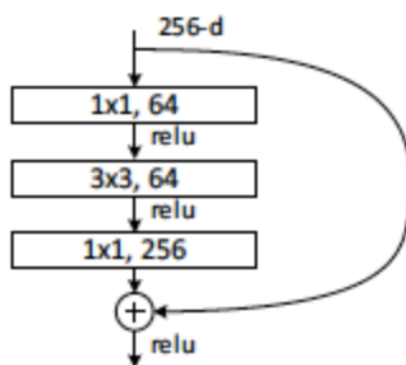


Figura 5.2: Esquema general de un bloque *Bottleneck* para una ResNet [44].

En este trabajo se experimenta de forma sistemática con varias profundidades de ResNet 3D: ResNet-18, ResNet-34, ResNet-101 y ResNet-152. La motivación es encontrar un equilibrio entre dos extremos: por un lado, redes demasiado simples que no capturen toda la complejidad de los patrones estructurales en las MRI cerebrales, y por otro, redes excesivamente profundas que resulten difíciles de entrenar con el tamaño de *dataset* disponible. Las variantes muy profundas, apoyadas en bloques *Bottleneck* (ResNet-101 y ResNet-152), muestran en la práctica un claro problema de infraajuste (*underfitting*): pese a su alta capacidad teórica, el modelo no es capaz de ajustar correctamente todos los parámetros con el número de sujetos disponible y obtiene resultados de validación inferiores a los de arquitecturas más pequeñas. Esto se refleja en pérdidas relativamente altas tanto en entrenamiento como en validación y en curvas de aprendizaje que no terminan de converger.

En el extremo opuesto, la ResNet-18 converge en muy pocas épocas debido a su menor profundidad y a que cuenta con menos parámetros, alcanzando alrededor de un 99 % de exactitud en entrenamiento, pero su rendimiento en validación se queda en torno al 69 %,

lo que indica un sobreajuste severo (*overfitting*). El modelo es capaz de memorizar bien los datos vistos, pero no generaliza cuando se le presentan cerebros nuevos. La ResNet-34 se sitúa entre ambos extremos y ofrece el mejor compromiso entre profundidad y capacidad de generalización: sus curvas de pérdida son más estables, la brecha entre entrenamiento y validación se reduce de forma significativa y las métricas de validación son consistentemente superiores al resto de variantes. Por estos motivos, y tras comparar cuantitativamente los resultados, se selecciona la ResNet-34 como arquitectura definitiva para los experimentos de clasificación que se describen en los capítulos posteriores. Finalmente, en la Figura 5.3 se muestra la estructura final de nuestro modelo de clasificación.

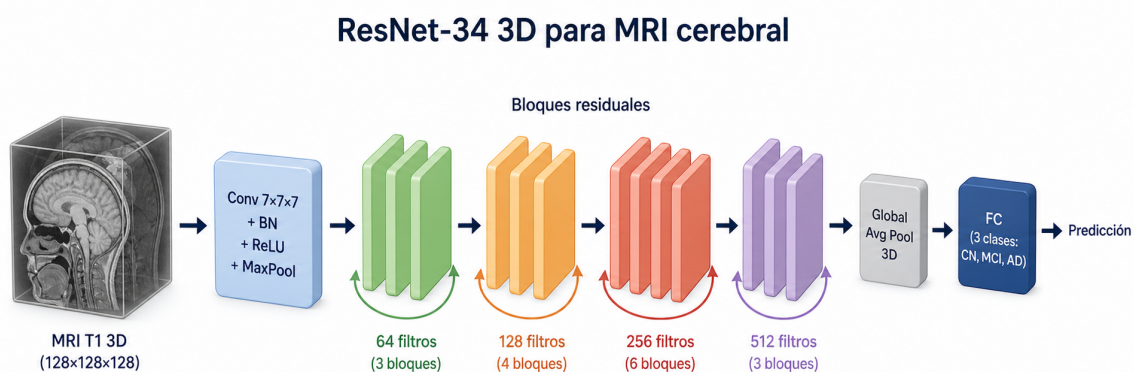


Figura 5.3: Esquema de la arquitectura ResNet-34 3D empleada para la clasificación de imágenes T1 cerebrales.

## 5.2 Entrenamiento del clasificador

El entrenamiento de los modelos se lleva a cabo en un clúster proporcionado por la Universidad Pontificia Comillas, basado en un nodo NVIDIA DGX con ocho GPUs H200. Esto supone una diferencia muy significativa respecto a utilizar una única GPU convencional: entrenamientos que en un ordenador personal requieren varias horas o incluso días por configuración se completan en del orden de minutos u pocas horas, y además es posible lanzar varios experimentos en paralelo. Gracias a este entorno de cómputo, resulta viable explorar distintas profundidades de ResNet, funciones de pérdida y combinaciones de hiperparámetros dentro del tiempo disponible para el TFG, algo que resulta muy difícil con recursos hardware domésticos.

El clasificador se entrena utilizando los volúmenes T1 preprocesados de dimensiones  $128 \times 128 \times 128$  y sus etiquetas clínicas organizadas en los conjuntos de entrenamiento, validación y prueba descritos en el capítulo previo. Se trabaja con un tamaño de *batch*

de 8 volúmenes, que resulta ser un compromiso razonable entre estabilidad del gradiente y memoria disponible en la GPU. Cada época de entrenamiento recorre por completo el conjunto de entrenamiento y, al finalizar, se evalúa el modelo sobre el conjunto de validación, registrando tanto la pérdida como la exactitud en ambos conjuntos.

En cuanto a la función de pérdida, se sigue un proceso incremental. Inicialmente se emplea la entropía cruzada estándar (*Cross-Entropy Loss*) para clasificación multiclase, que permite entrenar los primeros modelos y comparar arquitecturas. Sin embargo, el desbalance de clases del *dataset* (con más sujetos CN que MCI y AD) hace que el modelo tienda a favorecer la clase mayoritaria. Para corregirlo, se introducen pesos de clase en la entropía cruzada, asignando un peso mayor a CN y AD que a MCI, de forma que los errores en las clases minoritarias penalicen más durante el entrenamiento. Aun así, las dificultades para aprender correctamente los casos más ambiguos motivan un último cambio hacia una *Focal Loss* multiclase, que combina estos pesos de clase ( $\alpha = [1,5, 0,8, 2,5]$ ) con un factor de enfoque  $\gamma = 2,0$  para reducir el impacto de las muestras fáciles y concentrar el aprendizaje en las difíciles. Esta combinación mejora la sensibilidad en MCI y AD y estabiliza el entrenamiento frente al desbalance.

Como optimizador se utiliza AdamW, una variante de Adam que desacopla explícitamente el término de *weight decay* de la actualización de los parámetros, lo que proporciona una regularización más controlada y suele traducirse en mejores métricas de validación en comparación con Adam clásico. Se fija una tasa de aprendizaje (*learning rate*) inicial de  $3 \times 10^{-4}$  y un *weight decay* de  $10^{-2}$ . Sobre este optimizador se aplica un plan de adaptación automática de la tasa de aprendizaje mediante un *scheduler* de tipo *ReduceLROnPlateau*, que reduce la *tasa de aprendizaje* a la mitad cuando la pérdida de validación deja de mejorar durante varias épocas consecutivas (paciencia 5, con un valor mínimo de  $10^{-6}$ ). En la práctica, esto permite comenzar con pasos relativamente grandes y refinar progresivamente el ajuste a medida que el modelo se acerca a un mínimo.

Para controlar el sobreajuste se combinan varias técnicas de regularización. Además del *weight decay* ya mencionado, se emplea una capa de *dropout* con probabilidad 0,5 justo antes de la capa totalmente conectada final, reduciendo la coadaptación de las neuronas en las últimas capas densas. Se establece un máximo de 50 épocas de entrenamiento, pero se incorpora un mecanismo de *early stopping* con paciencia de 10 épocas sin mejora en la exactitud de validación. Cada vez que ésta supera el máximo anterior se guardan los pesos del modelo como “mejor *checkpoint*”; si tras 10 épocas consecutivas no se observa mejora, el entrenamiento se detiene automáticamente y se conserva el modelo con mejor rendimiento en validación. Esta combinación de pérdida focal con pesos de clase, optimización mediante AdamW, ajuste dinámico de la tasa de aprendizaje, *dropout* y *early stopping*, junto con el uso de un clúster DGX con GPUs H200, permite entrenar la ResNet-34 de

forma estable y obtener un clasificador capaz de generalizar razonablemente bien a sujetos no vistos, a pesar del tamaño limitado y el desbalance del conjunto de datos.

### 5.3 Métricas de evaluación

Para evaluar el rendimiento del clasificador se utilizan varias métricas estándar en problemas de clasificación multiclase, todas ellas derivadas de la matriz de confusión (verdaderos positivos, verdaderos negativos, falsos positivos y falsos negativos) calculada sobre el conjunto de prueba.

La métrica más sencilla es la exactitud, que mide el porcentaje de ejemplos correctamente clasificados sobre el total. Es útil como visión global, pero en presencia de clases desbalanceadas puede resultar engañosa, ya que un modelo que favorezca la clase mayoritaria puede obtener una exactitud alta sin comportarse bien en las clases minoritarias.

Para cada clase se definen, además, la *sensibilidad* (o *recall*) y la *especificidad*. La sensibilidad indica qué proporción de sujetos de una clase concreta el modelo identifica correctamente (verdaderos positivos entre todos los positivos reales), mientras que la especificidad mide qué proporción de sujetos que no pertenecen a esa clase se clasifican correctamente como negativos (verdaderos negativos entre todos los negativos reales). A partir de la matriz de confusión también se calcula la *precisión* de cada clase, que indica de todos los sujetos que el modelo predice como pertenecientes a esa clase qué porcentaje lo son realmente (verdaderos positivos entre todos los positivos predichos).

En términos formales, para una clase dada:

$$\text{Sensibilidad} = \frac{TP}{TP + FN}, \quad \text{Especificidad} = \frac{TN}{TN + FP}, \quad \text{Precisión} = \frac{TP}{TP + FP}.$$

La *F1-score* combina precisión y sensibilidad en una única métrica, mediante su media armónica, de forma que solo toma valores altos cuando ambas son elevadas. Esto resulta especialmente útil cuando se quiere equilibrar la capacidad del modelo para no dejar escapar casos positivos (alta sensibilidad) y para no producir demasiados falsos positivos (alta precisión), algo habitual en aplicaciones médicas. En este trabajo se consideran F1-scores por clase y también promedios globales: el promedio *macro*, que calcula la media de las métricas de cada clase dando el mismo peso a todas, y el promedio *weighted*, que pondera cada clase según su frecuencia en el conjunto de datos.

Para una clase dada, la F1 viene dada por:

$$F1 = 2 \cdot \frac{\text{Precisión} \cdot \text{Sensibilidad}}{\text{Precisión} + \text{Sensibilidad}},$$

y los promedios se calculan como:

$$F1_{\text{macro}} = \frac{1}{K} \sum_{k=1}^K F1_k, \quad F1_{\text{weighted}} = \sum_{k=1}^K w_k F1_k,$$

donde  $K$  es el número de clases y  $w_k$  es la proporción de muestras de la clase  $k$  en el conjunto de datos.

Finalmente, se emplea el área bajo la curva ROC (*AUC-ROC*) para evaluar la capacidad del modelo de separar cada clase frente al resto a lo largo de distintos umbrales de decisión. La curva ROC representa la relación entre la tasa de verdaderos positivos y la tasa de falsos positivos, y el AUC resume dicha curva en un solo valor entre 0,5 (comportamiento similar al azar) y 1 (separación perfecta). En escenarios multiclase, el AUC se calcula típicamente en modo *one-vs-rest* para cada clase y, de nuevo, se puede resumir mediante promedios macro o ponderados.

# Capítulo 6

## Metodología de generación

Debido a la necesidad de más datos para mejorar el performance del modelo, se plantea el uso de modelos generativos capaces de sintetizar volúmenes T1 cerebrales realistas, de tamaño  $128 \times 128 \times 128$ , que se parezcan lo máximo posible a los sujetos de ADNI. Para ello se ha diseñado una arquitectura final basada en *Wasserstein GAN* (WGAN) adaptada a datos 3D [45].

En lugar de entrenar un único modelo condicional explícito, en este trabajo se ha optado por entrenar un modelo independiente por clase (CN, MCI y AD), de modo que cada WGAN aprende la distribución específica de su grupo diagnóstico. El condicionamiento por clase se realiza, por tanto, a nivel de datos: el modelo que se entrena con sujetos AD generará cerebros sintéticos de tipo AD, mientras que el entrenado con CN generará cerebros de sujetos sanos.

### 6.1 Arquitectura del WGAN

En la Sección 3.2.3 de este trabajo se explican los fundamentos de los modelos generativos y, en particular, de las arquitecturas tipo GAN. Sin embargo, su aplicación práctica a volúmenes 3D requiere tomar decisiones específicas de diseño (tamaño y profundidad de la red, tipo de bloques residuales, función de pérdida, estrategia de entrenamiento, etc.) y ajustar varios hiperparámetros que afectan de forma directa a la estabilidad del entrenamiento y a la calidad visual de las imágenes generadas. Para llegar a la arquitectura final se fueron probando y refinando varios modelos previos, introduciendo cambios progresivos hasta encontrar una configuración que ofreciera un buen equilibrio entre estabilidad y calidad. Como paso inicial se implementó un modelo de tipo VAE-GAN, inspirado en el trabajo de Larsen et al. [46], que combina un autoencoder variacional con un discriminador adver-

sarial para mejorar las reconstrucciones. Este planteamiento facilita la generación inicial de volúmenes sintéticos tridimensionales; no obstante, adolece de limitaciones críticas: las imágenes carecen de la definición necesaria y el proceso de entrenamiento manifiesta una volatilidad extrema frente a los hiperparámetros, lo que impide su escalado eficiente. En la Figura 6.1 se muestra un ejemplo de salida de este modelo, donde se puede apreciar que la imagen es borrosa y no es válida si se quiere usar para los objetivos de este trabajo, ya que se requiere mucha más precisión por voxel. En consecuencia, aunque el VAE-GAN demuestra que es posible generar cerebros sintéticos a partir de los datos de ADNI, sus limitaciones en cuanto a nitidez y nivel de detalle motivaron la exploración de arquitecturas alternativas basadas en WGAN.

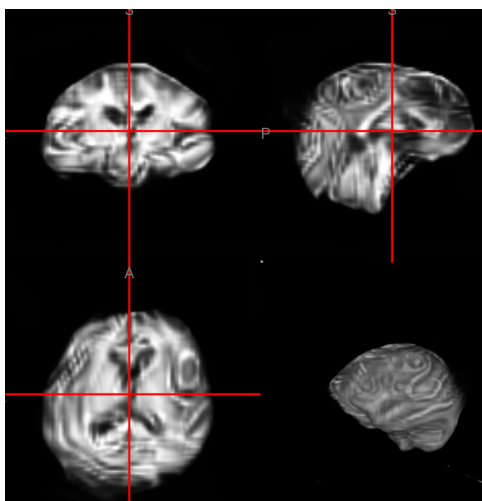


Figura 6.1: Ejemplo de salida del modelo inicial VAE-GAN.

A partir de estas limitaciones se decidió pasar a una familia de modelos basada en WGAN puro, sin la parte variacional, centrando el objetivo en aprender a generar directamente volúmenes realistas mediante un juego generador-discriminador más sencillo de controlar. Frente a los GAN clásicos, los WGAN utilizan una función de coste basada en la distancia de Wasserstein, lo que suele proporcionar un entrenamiento más estable y una mejor correlación entre la pérdida del discriminador y la calidad de las muestras generadas. Sobre esta base se probaron distintas versiones del generador y del discriminador (incluyendo cambios en la dimensión del espacio latente, en la forma de hacer el *upsampling* y en la normalización de las capas), hasta converger en la arquitectura WGAN 3D definitiva descrita en las subsecciones siguientes.

### 6.1.1 WGAN

En este trabajo se partió de un modelo WGAN 3D ya existente, concretamente la implementación pública *3dbraingen*, que genera volúmenes de resonancia magnética cerebral a

partir de un vector latente utilizando una arquitectura de tipo auto-encoding GAN [47]. Esta implementación, basada en el trabajo “Generation of 3D Brain MRI Using Auto-Encoding Generative Adversarial Network”, proporciona un punto de partida práctico, ya que incluye un encoder, un generador y un discriminador capaces de producir cerebros sintéticos plausibles a partir de datos de ADNI.

A partir de este modelo de referencia se fueron introduciendo cambios progresivos para adaptarlo a los objetivos concretos de este estudio. En primer lugar, se ajustaron la arquitectura y los tamaños de entrada/salida para trabajar con el preprocesado específico empleado en nuestro *pipeline* y con volúmenes T1 de tamaño  $128 \times 128 \times 128$ . Posteriormente, se modificó el espacio latente (reduciendo su dimensión y adaptando la forma de muestreo), se reorganizaron las redes de encoder y generador para integrarlas con los bloques residuales y las convoluciones transpuestas 3D, y se substituyó el discriminador original por un discriminador multi-escala con normalización espectral, más estable en nuestro contexto. Además, en lugar de entrenar un único modelo, se adoptó un esquema de varios WGAN independientes, uno por clase (CN, MCI y AD), para que cada modelo aprendiera mejor la distribución de su grupo diagnóstico. Tras esta serie de modificaciones se llegó a la arquitectura WGAN 3D final, que se describe con detalle en la siguiente subsección.

### 6.1.2 Arquitectura Final: WGAN crítico

En la Sección 3.2.3 se describen en detalle las diferentes partes de una GAN y su funcionamiento básico. Sobre esa base, la arquitectura final de nuestro modelo generativo está compuesta por cuatro bloques principales: un encoder, un generador, un discriminador multi-escala y un discriminador de código (*CodeDiscriminator*). De forma intuitiva, el encoder convierte un cerebro real en un vector numérico corto, el generador hace el camino inverso y fabrica un cerebro sintético a partir de ese vector, el discriminador multi-escala comprueba si un volumen “parece real” o generado, y el discriminador de código vigila que los vectores producidos por el encoder sigan una distribución latente razonable.

El **encoder** recibe como entrada un volumen T1 real de tamaño  $1 \times 128 \times 128 \times 128$ . Para procesarlo, aplica varias capas de convolución 3D con *stride* 2: cada vez que pasa por una de estas capas, la resolución espacial del volumen se reduce a la mitad ( $128 \rightarrow 64 \rightarrow 32 \rightarrow \dots$ ) mientras aumenta el número de canales. Cada bloque incluye además normalización y una activación no lineal, como *LeakyReLU*, elegida para favorecer un flujo de gradiente más estable y evitar que, si muchas activaciones resultan negativas al inicio del entrenamiento, neuronas completas dejen de aprender. Tras varias de estas capas, el volumen se ha convertido en un pequeño bloque 3D con muchos canales. Este bloque se

aplana y se pasa por una o varias capas totalmente conectadas, obteniendo finalmente un vector latente  $z \in \mathbb{R}^{256}$ . Ese vector de 256 números resume la forma y la intensidad del cerebro de entrada en un formato compacto que el generador puede usar.

El **generador** hace el proceso inverso: parte de un vector latente de dimensión 256 y reconstruye a partir de él un volumen T1 sintético de tamaño  $128 \times 128 \times 128$ . Primero, una capa totalmente conectada proyecta el vector  $z$  a un pequeño bloque 3D, por ejemplo de tamaño  $512 \times 4 \times 4 \times 4$ . A continuación, el generador aplica una secuencia de bloques **UpsampleBlock**: cada uno utiliza una convolución transpuesta 3D (que aumenta la resolución, por ejemplo  $4 \rightarrow 8 \rightarrow 16 \rightarrow 32 \rightarrow 64 \rightarrow 128$ ), seguida de normalización, activación y un bloque residual 3D (**ResBlock3D**) que ayuda a afinar detalles anatómicos sin destruir la estructura ya aprendida. Tras varios de estos bloques, se llega de nuevo a la resolución  $128^3$ . Una última convolución 3D reduce los canales a uno (una única imagen T1) y se aplica una función *tanh* para obtener intensidades normalizadas. En la práctica, el generador aprende a transformar puntos del espacio latente en cerebros 3D con apariencia realista.

El **crítico multi-escala** está formado por dos críticos WGAN que comparten la misma idea pero trabajan a distinta resolución: Un crítico recibe el volumen completo a resolución  $128^3$  y el otro recibe el mismo volumen reescalado a  $64^3$ .

Cada crítico es una red de convoluciones 3D con *stride 2* que va reduciendo la resolución del volumen ( $128 \rightarrow 64 \rightarrow 32 \rightarrow \dots$ ) hasta llegar a un solo número real como salida. Ese número indica cuánto se parece, según el modelo, el volumen de entrada a un cerebro real (valores altos) frente a uno generado (valores bajos). Al tener dos críticos, uno “mirando de lejos” ( $64^3$ ) y otro “de cerca” ( $128^3$ ), el generador se ve obligado a respetar tanto la forma global del cerebro como los detalles locales. Para que el entrenamiento sea más estable, las convoluciones de estos críticos se regularizan con técnicas como la normalización espectral, que controlan la magnitud de los pesos y evitan que el criticado se vuelva demasiado extremo.

El **discriminador de código** (*CodeDiscriminator*) trabaja en el espacio latente en lugar de en las imágenes. Es una red pequeña de capas totalmente conectadas que recibe vectores de dimensión 256 y debe decidir si: provienen del encoder aplicado a cerebros reales, o si han sido muestreados directamente de la distribución latente objetivo (por ejemplo, una distribución normal).

Durante el entrenamiento, el encoder se ajusta para “engañar” a este discriminador, es decir, para que sus códigos tengan la misma pinta que los vectores muestreados del prior, utilizando como señal de entrenamiento la salida escalar del propio discriminador de código (en lugar de una pérdida logarítmica clásica).

Esto hace que el espacio latente esté bien organizado: puntos cercanos en el espacio latente corresponden a cerebros parecidos, y el generador puede moverse suavemente por ese espacio para producir variaciones realistas.

En conjunto, la combinación encoder–generador permite ir y volver entre el espacio de imágenes 3D y un espacio latente compacto, mientras que el crítico multi-escala y el discriminador de código actúan como “jueces” que obligan a que tanto las imágenes generadas como los vectores latentes tengan un aspecto coherente con los datos reales. Gracias a ello, el modelo es capaz de generar volúmenes T1 sintéticos realistas por clase (CN, MCI, AD) y de disponer de una representación latente útil para tareas posteriores como la clasificación o el análisis de variaciones anatómicas.

## 6.2 Estrategia de entrenamiento

El entrenamiento del WGAN 3D se propone una combinación de tres objetivos: En primer lugar, que las imágenes generadas sean realistas (pérdida adversarial), en segundo lugar, las reconstrucciones a partir del encoder conserven bien la anatomía (pérdida de reconstrucción perceptual) y por último que el espacio latente tenga una distribución bien comportada (pérdida sobre el código). Todo ello se entrena utilizando la formulación Wasserstein con *gradient penalty*, que ofrece gradientes más estables que los GAN clásicos y una mejor correlación entre la pérdida del discriminador y la calidad de las muestras.

### 6.2.1 Funciones de pérdida

En cada iteración del entrenamiento se distinguen tres grupos de pérdidas: las asociadas al crítico multi-escala, las del discriminador de código y las correspondientes al bloque encoder–generador. Para formular las funciones de pérdida con las cuales se han optimizado los hiperparámetros, se usa la siguiente notación:  $x_{\text{real}}$  es un volumen real,  $x_{\text{rec}} = G(E(x_{\text{real}}))$  es su reconstrucción,  $x_{\text{rand}} = G(z_{\text{rand}})$  es una muestra generada a partir de un código gaussiano  $z_{\text{rand}} \sim \mathcal{N}(0, I)$ , y  $z_{\text{enc}} = E(x_{\text{real}})$  es el código latente producido por el encoder.

#### Pérdida del crítico multi-escala

El crítico multi-escala  $D$  recibe, para cada paso, volúmenes reales  $x_{\text{real}}$ , reconstrucciones  $x_{\text{rec}}$  y muestras puramente sintéticas  $x_{\text{rand}}$ . Para cada volumen, el discriminador produce dos salidas escalares: una a resolución completa  $128^3$ , que denotamos  $D_{\text{full}}(x)$ , y otra a resolución reducida  $64^3$ ,  $D_{\text{half}}(x)$ . La pérdida se construye siguiendo la formulación

WGAN, combinando en cada escala las contribuciones de reales, reconstrucciones y muestras aleatorias, y ponderando el término de los reales con un factor dos para equilibrar su influencia.

La pérdida adversarial del crítico multi-escala queda:

$$\mathcal{L}_D^{\text{adv}} = \left( D_{\text{full}}(x_{\text{rec}}) + D_{\text{full}}(x_{\text{rand}}) - 2 D_{\text{full}}(x_{\text{real}}) \right) + \left( D_{\text{half}}(x_{\text{rec}}) + D_{\text{half}}(x_{\text{rand}}) - 2 D_{\text{half}}(x_{\text{real}}) \right).$$

Sobre esta pérdida se añade un término de *gradient penalty* WGAN-GP, calculado sobre interpolaciones  $\hat{x}$  entre volúmenes reales y generados:

$$\mathcal{L}_{\text{GP}}^x = \lambda \mathbb{E}_{\hat{x}} \left[ \left( \|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1 \right)^2 \right],$$

donde  $\lambda = 5,0$  es el peso de la penalización y  $D(\hat{x}) = D_{\text{full}}(\hat{x}) + D_{\text{half}}(\hat{x})$ .

La pérdida total del crítico multi-escala es entonces

$$\mathcal{L}_D = \mathcal{L}_D^{\text{adv}} + \mathcal{L}_{\text{GP}}^x.$$

### Pérdida del discriminador de código

El discriminador de código  $CD$  opera en el espacio latente de dimensión 256 y recibe códigos procedentes del encoder  $z_{\text{enc}} = E(x_{\text{real}})$  y códigos gaussianos  $z_{\text{rand}} \sim \mathcal{N}(0, I)$ . Su objetivo es asignar valores altos a los códigos “reales” (gaussianos) y valores bajos a los producidos por el encoder, siguiendo también una formulación tipo WGAN.

La parte adversarial de la pérdida se escribe como

$$\mathcal{L}_{CD}^{\text{adv}} = \mathbb{E}[CD(z_{\text{enc}})] - \mathbb{E}[CD(z_{\text{rand}})].$$

De forma análoga al caso de las imágenes, se añade un término de *gradient penalty* en el espacio latente, calculado sobre interpolaciones  $\hat{z} = \alpha z_{\text{rand}} + (1 - \alpha) z_{\text{enc}}$ :

$$\mathcal{L}_{\text{GP}}^z = \lambda \mathbb{E}_{\hat{z}} \left[ \left( \|\nabla_{\hat{z}} CD(\hat{z})\|_2 - 1 \right)^2 \right].$$

La pérdida total del discriminador de código queda entonces

$$\mathcal{L}_{CD} = \mathcal{L}_{CD}^{\text{adv}} + \mathcal{L}_{\text{GP}}^z.$$

## Pérdida de reconstrucción perceptual

La pérdida de reconstrucción  $\mathcal{L}_{\text{rec}}$  se aplica sobre  $x_{\text{rec}} = G(E(x_{\text{real}}))$  y combina tres componentes: una pérdida L1 volumétrica, un término basado en SSIM y una *gradient loss* 3D.

La componente L1 volumétrica es

$$\mathcal{L}_{\text{L1}} = \mathbb{E} \left[ \|x_{\text{rec}} - x_{\text{real}}\|_1 \right].$$

Para la parte SSIM, se consideran tres cortes centrales (axial, coronal y sagital) de ambos volúmenes, que denotamos  $x_{\text{rec}}^{(d)}$  y  $x_{\text{real}}^{(d)}$  para cada plano  $d$ . Si  $\mathcal{S}(\cdot, \cdot)$  es la función SSIM, se define

$$\mathcal{L}_{\text{SSIM}} = \frac{1}{3} \sum_{d \in \{\text{axial}, \text{coronal}, \text{sagital}\}} \left( 1 - \mathcal{S}(x_{\text{rec}}^{(d)}, x_{\text{real}}^{(d)}) \right).$$

Por último, la *gradient loss* 3D compara los gradientes espaciales en las tres direcciones. Denotando por  $\nabla_x, \nabla_y, \nabla_z$  las diferencias finitas en cada eje, se tiene

$$\mathcal{L}_{\text{grad}} = \frac{1}{3} \left( \mathbb{E} \left[ \|\nabla_x x_{\text{rec}} - \nabla_x x_{\text{real}}\|_1 \right] + \mathbb{E} \left[ \|\nabla_y x_{\text{rec}} - \nabla_y x_{\text{real}}\|_1 \right] + \mathbb{E} \left[ \|\nabla_z x_{\text{rec}} - \nabla_z x_{\text{real}}\|_1 \right] \right).$$

La pérdida perceptual total utilizada en el entrenamiento se define como

$$\mathcal{L}_{\text{rec}} = \mathcal{L}_{\text{L1}} + 0,5 \mathcal{L}_{\text{SSIM}} + 0,3 \mathcal{L}_{\text{grad}}.$$

## Pérdidas de encoder y generador

El bloque encoder–generador se optimiza con la suma de dos términos: una pérdida de reconstrucción perceptual  $\mathcal{L}_{\text{rec}}$  y una pérdida adversarial  $\mathcal{L}_{\text{adv}}$ . La pérdida de reconstrucción se aplica sobre  $x_{\text{rec}} = G(E(x_{\text{real}}))$  y combina tres componentes que actúan de forma complementaria. En primer lugar, se calcula una pérdida L1 volumétrica entre  $x_{\text{rec}}$  y  $x_{\text{real}}$ , que penaliza diferencias de intensidad voxel a voxel. Esta componente se define como

$$\mathcal{L}_{\text{L1}} = \mathbb{E} \left[ \|x_{\text{rec}} - x_{\text{real}}\|_1 \right],$$

y fuerza a que las intensidades de la reconstrucción sean similares a las del volumen original en todo el cerebro.

En segundo lugar, se evalúa el índice de similitud estructural SSIM en los tres planos centrales (axial, coronal y sagital), obteniendo una medida de la similitud estructural

entre la reconstrucción y el original. Si denotamos por  $x_{\text{real}}^{(d)}$  y  $x_{\text{rec}}^{(d)}$  los cortes centrales en el plano  $d \in \{\text{axial, coronal, sagital}\}$ , y por  $\mathcal{S}(\cdot, \cdot)$  la función SSIM, el término asociado se define como

$$\mathcal{L}_{\text{SSIM}} = \frac{1}{3} \sum_{d \in \{\text{axial, coronal, sagital}\}} \left(1 - \mathcal{S}(x_{\text{rec}}^{(d)}, x_{\text{real}}^{(d)})\right).$$

De este modo, la pérdida no solo compara intensidades, sino también la estructura global de la imagen en los tres ejes principales.

En tercer lugar, se introduce una *gradient loss* 3D que compara los gradientes espaciales (derivadas en las tres direcciones) de ambos volúmenes. Denotando por  $\nabla_x, \nabla_y, \nabla_z$  las diferencias finitas a lo largo de cada eje, se define

$$\mathcal{L}_{\text{grad}} = \frac{1}{3} \left( \mathbb{E} \left[ \|\nabla_x x_{\text{rec}} - \nabla_x x_{\text{real}}\|_1 \right] + \mathbb{E} \left[ \|\nabla_y x_{\text{rec}} - \nabla_y x_{\text{real}}\|_1 \right] + \mathbb{E} \left[ \|\nabla_z x_{\text{rec}} - \nabla_z x_{\text{real}}\|_1 \right] \right).$$

Este término penaliza explícitamente las diferencias en los bordes y contornos, favoreciendo reconstrucciones con límites más nítidos y menos efecto de “plastilina”.

Las tres componentes se combinan con pesos fijos para formar la pérdida de reconstrucción perceptual total:

$$\mathcal{L}_{\text{rec}} = \mathcal{L}_{\text{L1}} + 0,5 \mathcal{L}_{\text{SSIM}} + 0,3 \mathcal{L}_{\text{grad}}.$$

Así,  $\mathcal{L}_{\text{rec}}$  captura simultáneamente similitud de intensidad, estructura global y nitidez local de los volúmenes generados.

La pérdida adversarial  $\mathcal{L}_{\text{adv}}$  se define a partir de las salidas del discriminador multi-escala y del discriminador de código. Por un lado, el encoder y el generador tratan de maximizar las salidas de  $D$  para las reconstrucciones  $x_{\text{rec}}$  y para las muestras aleatorias  $x_{\text{rand}} = G(z_{\text{rand}})$ , de manera que estos volúmenes se vuelvan indistinguibles de los reales a ojos de los críticos. Si denotamos por  $D_{\text{full}}(x)$  y  $D_{\text{half}}(x)$  las salidas de  $D$  a resolución completa y reducida, respectivamente, la contribución de estas dos escalas puede escribirse como

$$-2 \left( D_{\text{full}}(x_{\text{rec}}) + D_{\text{full}}(x_{\text{rand}}) + D_{\text{half}}(x_{\text{rec}}) + D_{\text{half}}(x_{\text{rand}}) \right),$$

donde el factor  $-2$  refleja que, desde el punto de vista de encoder y generador, se quiere invertir el papel del discriminador y hacer que sus salidas para imágenes generadas sean lo más altas posible.

Por otro lado, el encoder intenta reducir la capacidad de  $CD$  para diferenciar sus códigos de los muestreados directamente de la gaussiana. Esta presión se introduce mediante un término negativo con la salida del discriminador de código evaluado sobre

$$z_{\text{enc}} = E(x_{\text{real}}):$$

$$- CD(z_{\text{enc}}).$$

En conjunto, la pérdida adversarial que se deriva de estos dos módulos queda

$$\mathcal{L}_{\text{adv}} = -2 \left( D_{\text{full}}(x_{\text{rec}}) + D_{\text{full}}(x_{\text{rand}}) + D_{\text{half}}(x_{\text{rec}}) + D_{\text{half}}(x_{\text{rand}}) \right) - CD(z_{\text{enc}}).$$

Minimizar  $\mathcal{L}_{\text{adv}}$  empuja al encoder y al generador a producir volúmenes que el discriminador multi-escala considere realistas y códigos latentes que el discriminador de código no pueda distinguir de los gaussianos.

La pérdida total para encoder y generador combina la parte de reconstrucción y la parte adversarial mediante un peso de reconstrucción dependiente de la época:

$$\mathcal{L}_{E,G}(t) = w_{\text{rec}}(t) \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{adv}},$$

donde  $w_{\text{rec}}(t)$  se controla con un esquema de *warm-up* que comienza con un valor alto para forzar que el modelo aprenda primero a reconstruir bien la anatomía y va decreciendo de forma lineal hasta un valor final más bajo, permitiendo que la parte adversarial gane peso a medida que avanza el entrenamiento.

## 6.2.2 Peso de reconstrucción y número de pasos del discriminador

Durante el desarrollo del modelo se ha prestado especial atención a la relación entre el peso de reconstrucción, el número de actualizaciones del discriminador por cada paso de encoder–generador y la estabilidad general del entrenamiento.

### Evolución del peso de reconstrucción

El peso de reconstrucción  $w_{\text{rec}}(t)$  se controla con una función de *warm-up* que depende del número de época. En las primeras épocas el valor de  $w_{\text{rec}}$  es alto, lo que obliga al modelo a centrarse en aprender reconstrucciones anatómicamente correctas. A medida que avanza el entrenamiento, este peso decrece linealmente hasta alcanzar un valor final más bajo, que se mantiene constante hasta el final.

En la configuración inicial se utilizaron 80 épocas de *warm-up*, con un peso que disminuía desde 15.0 hasta 10.0, entrenando durante un total de 350 épocas con un tamaño de lote de 16. Tras analizar los resultados, se adoptó una configuración final en la que el *warm-up* se prolonga hasta 100 épocas, manteniendo el mismo rango de pesos, pero

aumentando el número total de épocas a 450 y reduciendo el tamaño de lote de 16 a 4. Aunque el entrenamiento con un lote de 16 era computacionalmente viable, al disponer de un conjunto de datos relativamente pequeño resultaba más recomendable trabajar con *batches* menores para evitar que demasiadas muestras quedaran agrupadas en cada actualización y mejorar así la capacidad de generalización del modelo. Este ajuste permite que el modelo dedique más tiempo a consolidar la estructura global de los cerebros antes de dar mayor protagonismo a la parte adversarial, y al mismo tiempo facilita el entrenamiento de arquitecturas más pesadas dentro de las limitaciones de memoria de la GPU.

### Número de pasos del discriminador por clase

El número de pasos del discriminador por cada actualización de encoder–generador (`n_critic`) se ha ajustado de forma específica para cada clase (CN, MCI y AD). En la configuración inicial, el modelo para la clase CN se entrenaba con cuatro pasos del discriminador por cada paso de encoder–generador, mientras que para MCI y AD se utilizaban dos. En la configuración final se redujo `n_critic_cn` a tres, manteniendo el valor de dos para MCI y AD. Esta pequeña variación permite que el generador asociado a la clase CN reciba relativamente más actualizaciones, acelerando su aprendizaje sin comprometer la estabilidad que proporciona un discriminador suficientemente entrenado.

### 6.2.3 Optimizadores y scheduler de la tasa de aprendizaje

Para cada uno de los cuatro bloques (encoder, generador, discriminador multi-escala y discriminador de código) se utiliza un optimizador Adam con parámetros  $\beta_1 = 0,0$  y  $\beta_2 = 0,9$ . Las tasas de aprendizaje se diferencian según el rol de cada red: el encoder y el generador emplean una tasa de  $1 \times 10^{-4}$ , mientras que el discriminador multi-escala y el discriminador de código utilizan una tasa de  $5 \times 10^{-5}$ . De este modo, las redes encargadas de producir las imágenes se adaptan algo más rápido, mientras que los discriminadores se actualizan de forma más conservadora.

Sobre cada optimizador se aplica un scheduler de tipo *CosineAnnealingLR*. Este scheduler reduce la tasa de aprendizaje de manera suave a lo largo de las épocas, empezando con pasos relativamente grandes que permiten explorar el espacio de parámetros y terminando con pasos más pequeños que favorecen la convergencia fina hacia una solución estable. El uso de este tipo de scheduler ha resultado especialmente relevante en entrenamientos largos, donde una tasa de aprendizaje fija tiende a producir oscilaciones en las últimas fases.

### 6.2.4 Técnicas de estabilización adicionales

Además de la formulación WGAN con *gradient penalty* y la normalización espectral en las convoluciones del discriminador, la estrategia de entrenamiento incluye varias decisiones orientadas a mejorar la estabilidad. La combinación de pérdida L1 volumétrica, SSIM en tres planos y gradient loss 3D en la pérdida de reconstrucción ayuda a reducir el típico problema de imágenes borrosas en modelos generativos, penalizando explícitamente la falta de nitidez en los bordes. El entrenamiento independiente de un modelo por clase (CN, MCI y AD) evita que un único discriminador deba abarcar distribuciones demasiado diferentes a la vez. Finalmente, el guardado periódico de *checkpoints* y de muestras intermedias permite monitorizar visualmente la evolución del modelo, detectar posibles inestabilidades (como colapso de modos) y ajustar los hiperparámetros según sea necesario.

## 6.3 Preprocesado y caché de volúmenes

El entrenamiento del WGAN se ha realizado, al igual que el modelo de clasificación, en el clúster de la universidad. Cada época podía tardar del orden de seis o siete minutos y, dado que para cada clase se entrena durante centenares de épocas (en la configuración final se emplean hasta 450 épocas por modelo), el tiempo de cómputo global resulta considerable. Una parte importante de este coste se debe a la carga y preprocesado de los volúmenes NIfTI en cada iteración, que incluye lectura desde disco, cambio de tamaño y normalización de intensidades.

Para mitigar este cuello de botella se implementó un sistema de caché de tensores. A partir del fichero de etiquetas, se recorre la lista de sujetos y, para cada uno de ellos, se localiza el archivo correspondiente en formato `.nii` o `.nii.gz`. El volumen se carga con `nibabel`, se convierte en un tensor de PyTorch de forma  $(1, X, Y, Z)$  y se envuelve en un sujeto de `TorchIO`. Sobre este sujeto se aplican las mismas transformaciones que se usarán durante el entrenamiento: un redimensionado a  $128 \times 128 \times 128$  y una normalización de intensidades al rango  $[-1, 1]$ , por ello, la salida del discriminador es una `tanh`. El tensor normalizado resultante se guarda en disco mediante `torch.save` en la carpeta de caché, utilizando el identificador del sujeto como nombre de fichero.

Durante el entrenamiento, el cargador de datos comprueba primero si existe el fichero `.pt` correspondiente en la caché. Si es así, se carga directamente el tensor preprocesado; solo en caso contrario se recurre al volumen NIfTI original y se aplica el pipeline de preprocesado. Gracias a este mecanismo, las épocas posteriores dejan de pagar el coste de E/S y transformación pesada, reduciendo de forma notable el tiempo por época y

permitiendo entrenar modelos generativos 3D complejos en un tiempo razonable dentro de los recursos del clúster.

## 6.4 Uso previsto de las imágenes sintéticas

El modelo WGAN 3D entrenado en este trabajo permite generar volúmenes T1 sintéticos específicos para cada clase (CN, MCI y AD). Estas imágenes no se conciben como un sustituto de los datos reales, sino como un recurso complementario que puede explotarse en tres direcciones principales: aumento de datos para el clasificador, análisis cualitativo mediante inspección visual y evaluación cuantitativa basada en métricas e impacto sobre el rendimiento de clasificación. La literatura previa ha mostrado que, en otros contextos de imagen médica, el uso de imágenes sintéticas generadas con GAN puede mejorar de forma apreciable la sensibilidad y la especificidad de los clasificadores al enriquecer el conjunto de entrenamiento y cubrir mejor la variabilidad de la población.

### 6.4.1 Aumento de datos para el clasificador

El uso más directo de las imágenes sintéticas es como mecanismo de *aumento de datos* para el clasificador 3D entrenado en capítulos anteriores. La idea es generar, para cada clase, un conjunto adicional de volúmenes T1 que sigan la distribución aprendida por el WGAN correspondiente (CN, MCI o AD) y mezclarlos con los volúmenes reales durante el entrenamiento del clasificador. Esta estrategia es especialmente interesante en escenarios donde algunas clases están infrarrepresentadas (por ejemplo, MCI o AD frente a CN), ya que permite equilibrar mejor el número de ejemplos por grupo sin necesidad de nuevas adquisiciones reales. De forma análoga a trabajos previos en otras modalidades médicas [48], se plantea comparar el rendimiento del clasificador entrenado solo con datos reales frente a la versión entrenada con datos reales más sintéticos, midiendo si se obtiene una mejora en métricas como la exactitud, la sensibilidad o el área bajo la curva ROC.

### 6.4.2 Análisis cualitativo y exploración visual

Además del uso cuantitativo, las imágenes generadas sirven como herramienta para un análisis cualitativo de la representación aprendida por el modelo. Al muestrear distintos vectores latentes y visualizar los volúmenes resultantes en los tres planos ortogonales, es posible inspeccionar si el WGAN captura patrones anatómicos característicos de cada grupo diagnóstico, como el patrón de atrofia en AD o la preservación estructural en CN. También se pueden comparar reconstrucciones  $x_{\text{rec}} = G(E(x_{\text{real}}))$  con sus correspondientes

volúmenes reales, evaluando visualmente en qué regiones la reconstrucción es más precisa y dónde aparecen artefactos o pérdida de detalle. Este tipo de inspección permite detectar fallos que no siempre se reflejan de forma clara en las métricas numéricas, y sirve como validación cualitativa de que el modelo no está aprendiendo estructuras anatómicamente inverosímiles.

### 6.4.3 Evaluación cuantitativa e impacto en el clasificador

Por último, las imágenes sintéticas se utilizan para una evaluación cuantitativa del modelo generativo y de su utilidad práctica. A nivel de imagen, se pueden calcular métricas de similitud entre reconstrucciones y volúmenes reales (por ejemplo, SSIM promedio en los planos ortogonales o medidas basadas en la pérdida perceptual definida en la sección de entrenamiento), así como estadísticas básicas de intensidad para comprobar que los cerebros generados respetan los rangos y distribuciones observados en los datos de ADNI [28], [49]. Además, se emplean métricas específicas de modelos generativos, como la *Fréchet Inception Distance* (FID) y, en su caso, el *Inception Score* (IS), que permiten cuantificar hasta qué punto la distribución de las imágenes sintéticas se aproxima a la de las imágenes reales [25], [26].

A nivel funcional, la evaluación más relevante consiste en medir el impacto que tiene la inclusión de imágenes sintéticas en el rendimiento del clasificador: se entrena un modelo de referencia exclusivamente con datos reales y se compara con un modelo entrenado con la misma arquitectura pero usando además muestras generadas por los WGAN de cada clase. La diferencia en métricas de clasificación (precisión, sensibilidad, especificidad, AUC) ofrece una medida indirecta de la calidad y utilidad de las imágenes sintéticas, siguiendo una estrategia similar a la utilizada en otros trabajos donde el beneficio de la síntesis se evalúa a través de tareas de diagnóstico automatizado [48].

# Capítulo 7

## Resultados

### 7.1 Resultados del modelo de clasificación

En esta sección se presentan los resultados del modelo de clasificación 3D entrenado sobre los volúmenes T1 de la base de datos ADNI, considerando tres clases diagnósticas: sujetos cognitivamente normales (CN), pacientes con deterioro cognitivo leve (MCI) y pacientes con enfermedad de Alzheimer (AD). Tras el preprocesado descrito en capítulos anteriores, el conjunto total disponible para este trabajo queda formado por 458 volúmenes CN, 231 volúmenes MCI y 193 volúmenes AD, que se dividen en subconjuntos de entrenamiento, validación y test manteniendo el equilibrio entre clases en la medida de lo posible.

A partir de esta partición se evalúan varias configuraciones del clasificador 3D, modificando tanto la arquitectura (por ejemplo, profundidad y número de filtros) como hiperparámetros de entrenamiento (tasa de aprendizaje, *dropout*) y la estrategia de enriquecimiento de datos, incluyendo o no técnicas de *data augmentation* clásica. El objetivo de esta comparación es identificar qué combinación de arquitectura y configuración de entrenamiento ofrece el mejor compromiso entre exactitud global, equilibrio entre clases y estabilidad de las métricas en validación y test.

Los resultados obtenidos se sitúan dentro del rango reportado por trabajos previos que abordan la clasificación CN/MCI/AD con datos de ADNI, donde se comunican exactitudes típicas en torno al 70–80 % para la clasificación en tres clases utilizando T1 cerebral, con conjuntos de entre unos cientos y algunos miles de sujetos [50], [51], [52]. Por ejemplo, Bermúdez et al. generan MRI sintéticas a partir de ADNI y reportan una exactitud del 74.5 % en la clasificación CN/MCI/AD sobre imágenes T1 3T, superando a modelos puramente discriminativos entrenados con los mismos datos [50]. De forma similar, en 3D-MobiBrainNet se emplean 221 sujetos con Alzheimer, 477 MCI y 284 CN de ADNI para

entrenar una red 3D multicategoría, obteniendo *accuracies* en torno al 80 % según la tarea considerada [51]. Además, estudios de revisión sobre clasificación multiclase HC/MCI/AD con MRI estructural recogen rangos de exactitud entre aproximadamente el 59 % y el 77 %, en función del modelo y del tamaño muestral empleado [52]. En los siguientes apartados se detallan las métricas de validación y test para cada configuración evaluada y se compara el rendimiento del mejor modelo con estas referencias de la literatura.

Como punto de partida se entrena una arquitectura ResNet18 tridimensional adaptada a los volúmenes T1 preprocesados de ADNI, utilizando las mismas particiones de entrenamiento, validación y test descritas anteriormente (636 volúmenes en entrenamiento, 113 en validación y 133 en test). La Figura 7.1 recoge la evolución de la precisión tanto en entrenamiento como en validación a lo largo de 50 épocas.

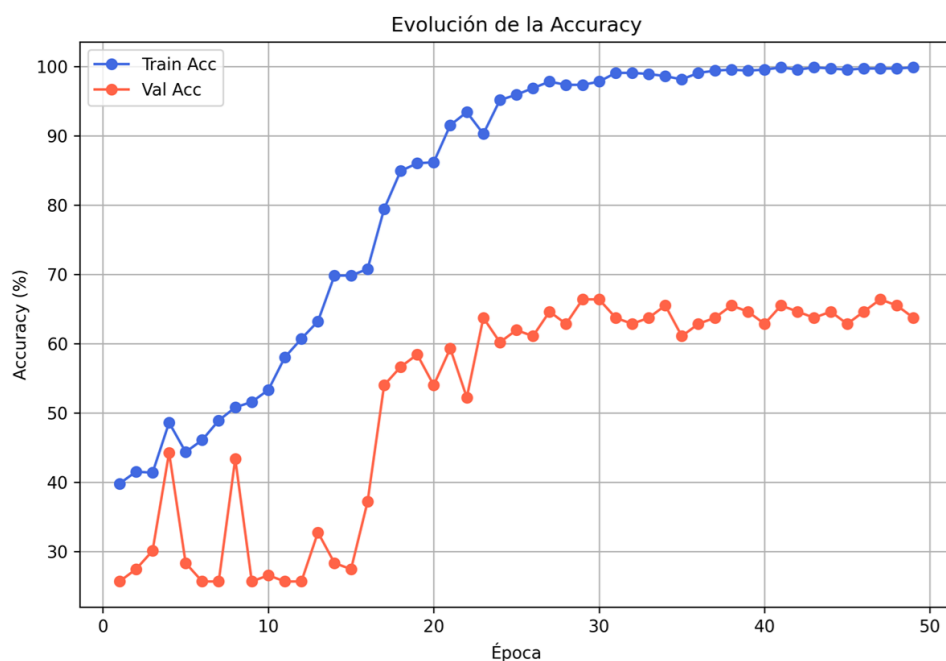


Figura 7.1: Curva de exactitud para el modelo ResNet18 en los conjuntos de entrenamiento y validación.

En la gráfica se aprecia cómo la *exactitud de validación* presenta una variabilidad elevada en las primeras épocas, este comportamiento se explica porque, en la fase inicial del entrenamiento, el modelo explora un espacio muy amplio de soluciones: los pesos se actualizan con pasos relativamente grandes y la pérdida en validación oscila de forma acusada hasta que la red empieza a especializarse en patrones consistentes.

La combinación de una exactitud de entrenamiento cercana al 100 % con una exactitud de validación que deja de mejorar y permanece significativamente separada indica un sobreajuste pronunciado del modelo. En otras palabras, la ResNet18 termina memorizando patrones específicos del subconjunto de entrenamiento, incluyendo ruido y posibles sesgos

de la muestra, pero no consigue capturar de forma robusta las diferencias entre CN, MCI y AD que se manifiestan en el conjunto de validación.

Para cuantificar el rendimiento del modelo se selecciona el mejor punto según la precisión de validación, correspondiente a la época 27, con una *exactitud de validación* del 64.60%. Evaluado sobre el conjunto de test, este punto alcanza una *pérdida de test* de 1.0166 y una *exactitud de test* del 69.92%, todavía baja comparada con otras investigaciones de la clasificación CN/MCI/AD con MRI estructural (aproximadamente 59–77%).

Las métricas de test por clase se resumen en la Tabla 7.1.

Clase	Precision	Recall	F1.
CN	0.78	0.71	0.74
MCI	0.52	0.66	0.58
AD	<b>0.81</b>	<b>0.72</b>	<b>0.76</b>

Tabla 7.1: Métricas de test para el modelo ResNet18.

Real/Pred	CN	MCI	AD
CN	<b>49</b>	17	3
MCI	10	<b>23</b>	2
AD	4	4	<b>21</b>

Tabla 7.2: Matriz de confusión en test para el modelo ResNet18.

El informe de clasificación muestra que el modelo resulta más fiable para la clase AD y, en menor medida, para CN, mientras que la clase MCI presenta una precisión claramente inferior y concentra la mayor parte de los errores, en línea con las confusiones observadas en la matriz de confusión entre CN y MCI. En conjunto, estos resultados confirman un sobreajuste acusado: el modelo logra prácticamente memorizar el conjunto de entrenamiento, pero solo alcanza un rendimiento moderado en validación y test, especialmente en la clase MCI. No obstante, este primer experimento resulta clave para detectar estas limitaciones y motiva el diseño de arquitecturas posteriores más profundas y mejor regularizadas, mejor adaptadas a los datos disponibles.

Este comportamiento sugiere que la ResNet18 empleada en este primer experimento tiene una capacidad excesiva para el tamaño efectivo del conjunto de entrenamiento y el desequilibrio entre clases. Aun así, este modelo resulta clave como punto de partida, ya que permite detectar estas limitaciones e identificar la necesidad de introducir mayor regularización y ajustar la arquitectura, lo que motiva la exploración de configuraciones posteriores basadas en redes más profundas y mejor adaptadas a los datos disponibles.

Además de la ResNet18, se evalúa una arquitectura considerablemente más profunda basada en ResNet152, con el objetivo de comprobar si un aumento de capacidad puede traducirse en una mejora del rendimiento. Sin embargo, el número de parámetros de esta red resulta desproporcionado para el tamaño efectivo del conjunto de entrenamiento, de modo que el optimizador no logra ajustar correctamente todos los pesos y el modelo queda

atrapado en una solución pobre. En concreto, aunque el mejor punto en validación alcanza una *exactitud de validación* del 67.50 %, la evaluación sobre el conjunto de test arroja una *exactitud de test* de solo el 51.88 %. Dado que la clase mayoritaria en el conjunto de datos es CN, el modelo termina colapsando hacia esta clase: en test predice prácticamente siempre CN, ignorando AD y MCI, como se aprecia en las métricas de la Tabla 7.3 y en la matriz de confusión correspondiente. En vista de estos resultados claramente inferiores a los obtenidos con la ResNet18, y del riesgo añadido de sobreajuste que implica una arquitectura tan profunda para un conjunto de datos limitado y desequilibrado, se descarta la ResNet152 y se opta por explorar modelos de capacidad intermedia y estrategias de regularización más adecuadas para el escenario planteado. Este comportamiento es coherente con lo que cabría esperar desde el punto de vista teórico: cuando la capacidad del modelo crece mucho más rápido que la cantidad de información disponible, la red dispone de suficientes grados de libertad como para minimizar parcialmente la pérdida apoyándose casi exclusivamente en la clase dominante, sin llegar a aprender representaciones discriminativas útiles para las clases minoritarias, lo que se traduce precisamente en el colapso hacia CN observado en test.

Clase	Precision	Recall	F1-score
CN	0.52	1.00	0.68
MCI	0.00	0.00	0.00
AD	0.00	0.00	0.00

Tabla 7.3: Métricas de test para el modelo ResNet152.

<b>Real/Pred</b>	CN	MCI	AD
CN	<b>69</b>	0	0
MCI	<b>35</b>	0	0
AD	<b>29</b>	0	0

Tabla 7.4: Matriz de confusión en test para el modelo ResNet152.

Finalmente, a la vista de los resultados obtenidos con la ResNet18 y la ResNet152, se opta por emplear una arquitectura intermedia basada en ResNet34, que ofrece un compromiso más adecuado entre capacidad del modelo y tamaño del conjunto de datos. Con la misma partición de 636 volúmenes para entrenamiento, 113 para validación y 133 para test, y un tamaño de *batch* de 4, la Figura 7.2 muestra la evolución de la precisión durante el entrenamiento. En este caso, la *exactitud de validación* alcanza valores cercanos al 70–75 %, con una separación más contenida entre las curvas de entrenamiento y validación que en la ResNet18. Esto indica un sobreajuste más controlado, donde la curva de entrenamiento continúa por encima de la de validación, aunque la diferencia es menor que en los modelos anteriores.

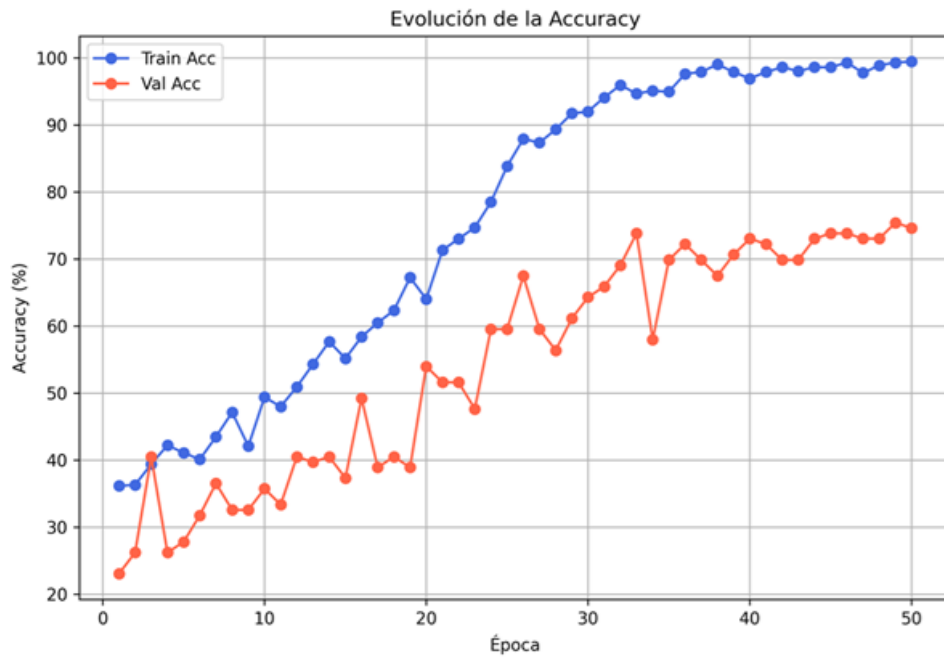


Figura 7.2: Curva de exactitud para el modelo ResNet34 en los conjuntos de entrenamiento y validación.

A partir de esta curva se selecciona como modelo final la época 30, correspondiente a una *exactitud de validación* del 73.45%. Evaluado sobre el conjunto de test, este punto alcanza una *exactitud de test* del 73.68%, mejorando claramente los resultados de la ResNet18. Las métricas de test por clase y la matriz de confusión se recogen en la Tabla 7.5 y la Tabla 7.6; en ellas se aprecia un comportamiento equilibrado para AD y CN (F1 en torno a 0.77–0.78) y una mejora respecto a la ResNet18 en la clase MCI (precisión 0.80, aunque con una sensibilidad más moderada), lo que sugiere que la ResNet34 captura mejor las diferencias sutiles asociadas a esta categoría intermedia.

Clase	Precision	Recall	F1
CN	0.72	<b>0.84</b>	<b>0.78</b>
MCI	<b>0.80</b>	0.46	0.58
AD	0.73	0.83	0.77

Tabla 7.5: Métricas de test para el modelo ResNet34.

<b>Real/Pred</b>	CN	MCI	AD
CN	<b>58</b>	4	7
MCI	<b>17</b>	16	2
AD	5	0	<b>24</b>

Tabla 7.6: Matriz de confusión en test para el modelo ResNet34.

El modelo final basado en ResNet34 presenta un rendimiento competitivo frente a otros trabajos similares, especialmente si se tiene en cuenta el tamaño relativamente reducido del conjunto de datos empleado. Mientras que numerosos estudios alcanzan precisiones similares o ligeramente superiores utilizando varios miles de volúmenes T1 o combinan-

do información multimodal (neuroimagen, pruebas cognitivas, biomarcadores), el modelo propuesto logra alrededor del 74% de exactitud en test con menos de 900 sujetos, manteniendo además un comportamiento equilibrado entre clases y una sensibilidad razonable para MCI.

Esta relación entre precisión conseguida y número de imágenes disponibles indica que la ResNet34 propuesta constituye una solución especialmente eficiente en escenarios con recursos limitados, donde no es factible entrenar redes de gran tamaño sobre cohortes muy extensas. En consecuencia, se considera un modelo robusto y potente en el contexto de la clasificación CN/MCI/AD con MRI estructural, ya que ofrece un rendimiento comparable al de enfoques más complejos, pero con un presupuesto de datos y de parámetros significativamente menor.

## 7.2 Clasificación binaria CN vs. Enfermedad

Tras completar la tarea de clasificación en tres clases (CN, MCI y AD), se pone de manifiesto que la principal dificultad clínica y algorítmica no reside tanto en distinguir a los pacientes ya diagnosticados de Alzheimer, sino en detectar de forma temprana la fase previa de deterioro cognitivo leve (MCI), que a menudo actúa como estado intermedio hacia la demencia. Dado que muchos trabajos previos señalan también la complejidad de separar MCI de CN y de AD en esquemas multiclase, se plantea una reformulación del problema orientada a simplificar la decisión y focalizar el modelo en la detección de sujetos con posible patología.

Con este objetivo, se define una nueva tarea de clasificación binaria en la que el conjunto original se reetiqueta en dos grupos: sujetos cognitivamente normales (CN) frente a un grupo denominado *Enfermedad*, que agrupa tanto a los pacientes con MCI como a los pacientes con diagnóstico de Alzheimer (AD). De este modo, el modelo pasa a aprender una frontera de decisión entre cerebros sanos y cerebros con signos de neurodegeneración, reservando para trabajos futuros una caracterización más fina de las distintas fases dentro del grupo patológico, pero permitiendo ya estudiar hasta qué punto las representaciones aprendidas en la clasificación CN/MCI/AD son capaces de generalizar a este escenario binario CN vs. Enfermedad.

Para llevar a cabo este cambio de planteamiento es necesario modificar varios elementos de la estructura final. En primer lugar, se actualiza el fichero `csv` de anotaciones para reflejar la nueva codificación binaria, asignando la etiqueta 0 a los sujetos del grupo *Enfermedad* (MCI+AD) y la etiqueta 1 a los sujetos CN. En segundo lugar, se adapta la capa de salida de la red, reduciendo el número de neuronas de tres a dos (o a una

con activación sigmoïdal, según la variante empleada), de modo que la ResNet18 quede específicamente configurada para un problema de dos clases. Finalmente, se sustituye la función de pérdida multiclase utilizada en los modelos anteriores por una entropía cruzada binaria coherente con este nuevo esquema de etiquetas.

### 7.2.1 Resultados con ResNet18 binaria

Para esta tarea binaria se emplea una ResNet18 tridimensional, manteniendo la misma partición de datos que en el escenario de tres clases (636 volúmenes para entrenamiento, 113 para validación y 133 para test), si bien con la nueva agrupación CN vs. Enfermedad (458 CN y 424 casos patológicos). El proceso de entrenamiento se resume en la Figura 7.3, donde se observa que la *exactitud de validación* alcanza valores en torno al 75–85 %, lo que indica un equilibrio razonable entre ajuste y capacidad de generalización.

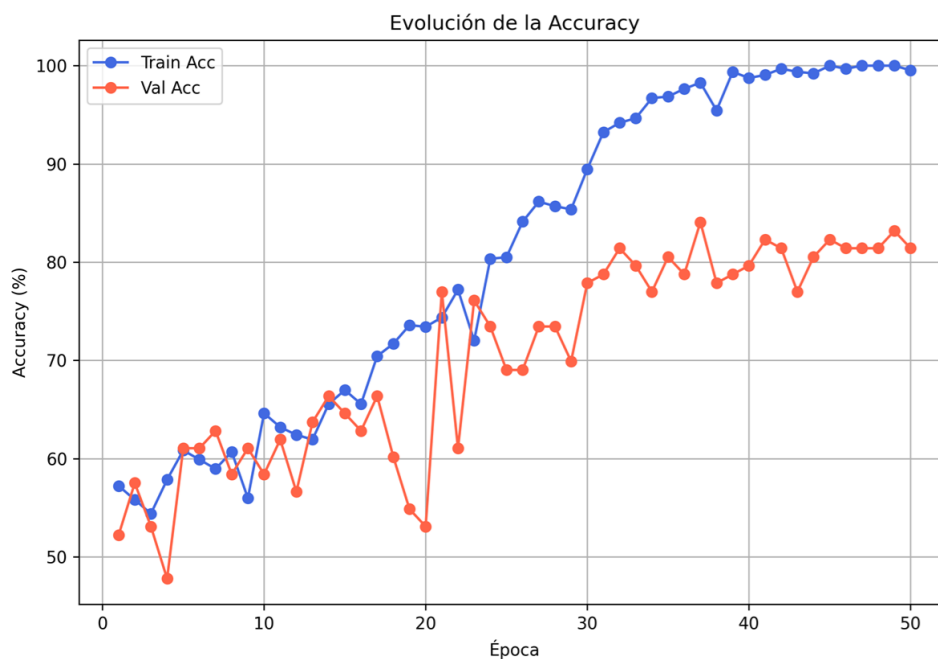


Figura 7.3: Curva de exactitud para la ResNet18 en la tarea binaria CN vs. Enfermedad.

El mejor modelo en validación se alcanza en la época 37, con una exactitud de validación del 84.07 %. Pese que sobre la época 20 ya alcanza valores óptimos de precisión, esta continúa incrementándose ligeramente hasta estabilizarse en torno al 84 %. Esto indica que, en esa fase del entrenamiento, el modelo ya no reduce de manera sustancial el error medio, pero sí ajusta la frontera de decisión de forma más fina y consigue clasificar correctamente algunos casos adicionales en validación, por lo que resulta razonable seleccionar la época 37 como punto óptimo. Evaluado sobre el conjunto de test, este modelo obtiene una test loss de 0.5075 y una test exactitud del 83.46 %, como se resume en la

Tabla 7.7. Estos valores suponen una mejora notable respecto al escenario multiclase, lo que confirma que la agrupación de MCI y AD en un único grupo patológico facilita la tarea al modelo y refuerza la idea de utilizar esta configuración como paso intermedio hacia esquemas de diagnóstico más finos.

Clase	Precision	Recall	F1
Enfermedad	<b>0.84</b>	0.81	0.83
CN	0.83	<b>0.86</b>	<b>0.84</b>

Tabla 7.7: Métricas de test para la ResNet18 en la tarea binaria CN vs. Enfermedad.

Real/Pred	Enfermedad	CN
Enfermedad	<b>52</b>	12
CN	10	<b>59</b>

Tabla 7.8: Matriz de confusión en test para la ResNet18 en la tarea binaria CN vs. Enfermedad.

Las métricas por clase muestran un comportamiento muy equilibrado entre CN y Enfermedad, con precisiones y sensibilidades en torno al 0.83 y un *F1-score* similar en ambas categorías. En conjunto, estos resultados indican que la ResNet18 binaria es capaz de discriminar de forma robusta entre sujetos sanos y sujetos con signos de patología neurodegenerativa a partir de volúmenes T1, constituyendo una base sólida sobre la que explorar estrategias posteriores de detección temprana específicas para MCI.

Además, el hecho de que el modelo mantenga un rendimiento similar en ambas clases sugiere que no se limita a explotar el ligero desequilibrio del conjunto de datos, sino que realmente aprende patrones estructurales asociados a la presencia de enfermedad. De este modo, la reformulación binaria CN vs. Enfermedad cumple el objetivo planteado: demuestra que, incluso con un número moderado de sujetos, es posible entrenar una red relativamente compacta (ResNet18) que identifica de forma fiable cerebros con alteraciones compatibles con MCI/AD, validando así la hipótesis de que la detección robusta del «estado patológico» puede servir como paso previo necesario para desarrollar modelos más finos orientados específicamente a diferenciar y caracterizar la fase de MCI.

### 7.3 Resultados del modelo de generación

En esta sección se presentan los resultados de los modelos generativos 3D entrenados sobre los volúmenes T1 de la base de datos ADNI, con el objetivo de sintetizar cerebros representativos de las tres clases diagnósticas consideradas: sujetos cognitivamente normales (CN), pacientes con deterioro cognitivo leve (MCI) y pacientes con enfermedad de Alzheimer (AD). En lugar de entrenar un único generador condicionado, se opta por entrenar tres modelos independientes, uno por cada clase, de manera que cada uno aprendiese de forma específica la distribución anatómica correspondiente a su grupo diagnóstico.

### 7.3.1 Modelo final de generación

El modelo generativo finalmente empleado para esta parte del trabajo consiste en una arquitectura tipo WGAN 3D entrenada de forma independiente para cada clase diagnóstica. A nivel de configuración, cada modelo se entrena con un tamaño de *batch* de 4, un espacio latente de dimensión 256 y un total de 450 épocas. El *batch size* de 4 se fija como compromiso entre estabilidad numérica y las limitaciones de memoria derivadas de trabajar con volúmenes 3D de  $128^3$  vóxeles, mientras que la dimensión latente 256 resulta suficiente para capturar la variabilidad anatómica observada en las MRI sin introducir un espacio de búsqueda excesivamente grande que dificulte la convergencia. El número de 450 épocas se establece empíricamente a partir de las curvas de pérdida y de la inspección visual de las muestras generadas, observándose que a partir de ese punto la calidad percibida apenas mejora y, en cambio, aumenta el riesgo de sobreajuste del discriminador. El optimizador utilizado fue Adam con tasa de aprendizaje diferenciada para generador/-codificador y discriminadores, empleando además un esquema de decaimiento cosenoidal de la tasa de aprendizaje. Asimismo, el número de actualizaciones del discriminador por cada actualización del generador se ajustó de forma específica para cada clase; en el caso de MCI, se utilizaron dos pasos del crítico por iteración de generador, mientras que para CN se emplearon tres y para AD dos, buscando en cada caso un equilibrio adecuado entre estabilidad adversaria y velocidad de convergencia.

La Figura 7.4 muestra la evolución cualitativa de las muestras generadas a lo largo del entrenamiento para distintas épocas representativas: 5, 20, 50, 100, 185, 330 y 450. Esta secuencia resulta especialmente útil para analizar cómo el modelo pasa de una representación inicial prácticamente amorfa a una estructura cerebral progresivamente más reconocible y anatómicamente coherente.

En la época 5, el modelo todavía se encuentra en una fase muy temprana de aprendizaje y las imágenes generadas presentan una estructura muy rudimentaria. Se intuye ya una cierta concentración de intensidad en la región central, pero la morfología cerebral aún no está bien definida: los bordes son difusos, existe bastante ruido y la forma global del volumen sigue siendo poco realista. En esta etapa, el generador apenas ha aprendido una aproximación gruesa a la distribución de intensidades de los cerebros reales.

En la época 20 se aprecia una mejora clara respecto a la fase inicial. La silueta cerebral comienza a consolidarse y ya se distinguen mejor los planos principales, especialmente en las vistas coronal y axial. Sin embargo, la imagen sigue mostrando una textura muy borrosa y una definición limitada de las estructuras internas. El modelo ha comenzado a capturar la geometría general del encéfalo, pero todavía no representa con suficiente precisión ni los ventrículos ni la separación entre tejidos.

A partir de la época 50 las imágenes generadas muestran ya una estructura cerebral claramente identificable. El contorno externo aparece mejor delimitado, la forma global es estable y comienzan a apreciarse zonas internas diferenciadas, si bien persisten áreas borrosas y cierto nivel de ruido residual.

En la época 100 se observa una mejora adicional en la organización espacial del volumen. Las estructuras internas se distribuyen de manera más equilibrada, las diferentes vistas presentan una correspondencia más coherente entre sí y la morfología global resulta más compacta y regular que en las épocas anteriores, lo que indica que el modelo ha captado gran parte de la variabilidad media de la clase MCI.

En torno a la época 185 el modelo alcanza uno de sus mejores compromisos entre forma global y detalle. La silueta externa se muestra más limpia, el fondo presenta menos artefactos y la disposición interna de las estructuras es más ordenada, de modo que el volumen adquiere una configuración anatómica coherente, aunque todavía con una nitidez inferior a la de una MRI real.

En la época 330 la forma cerebral se mantiene estable y anatómicamente razonable, pero se aprecian ciertos patrones de textura repetitivos y una apariencia algo suavizada en determinadas regiones. Esto sugiere que, a partir de este punto, las mejoras cualitativas respecto a épocas intermedias son más graduales, y que el modelo ha consolidado principalmente la macroestructura mientras resulta más difícil seguir refinando el detalle fino.

En la época 450, correspondiente al final del entrenamiento, la anatomía global se encuentra bien establecida, las tres vistas principales son coherentes entre sí y el volumen presenta una morfología compatible con la de un cerebro obtenido mediante MRI. Sin embargo, las diferencias respecto a épocas inmediatamente anteriores son relativamente pequeñas, lo que indica que el modelo llevaba ya varias decenas de épocas próximo a un régimen de convergencia estable.

En conjunto, la evolución visual observada confirma que el modelo aprende de forma progresiva la distribución anatómica asociada a la clase CN: las primeras fases se centran en capturar la forma global del volumen, mientras que las épocas intermedias y finales se orientan a mejorar la organización interna y a estabilizar la apariencia general de las muestras. Aunque las imágenes generadas mantienen cierto grado de suavizado y artefactos propios del entrenamiento adversario en 3D, los resultados obtenidos muestran una estructura suficientemente coherente como para considerar estas muestras en análisis cualitativos y en experimentos de aumento de datos.

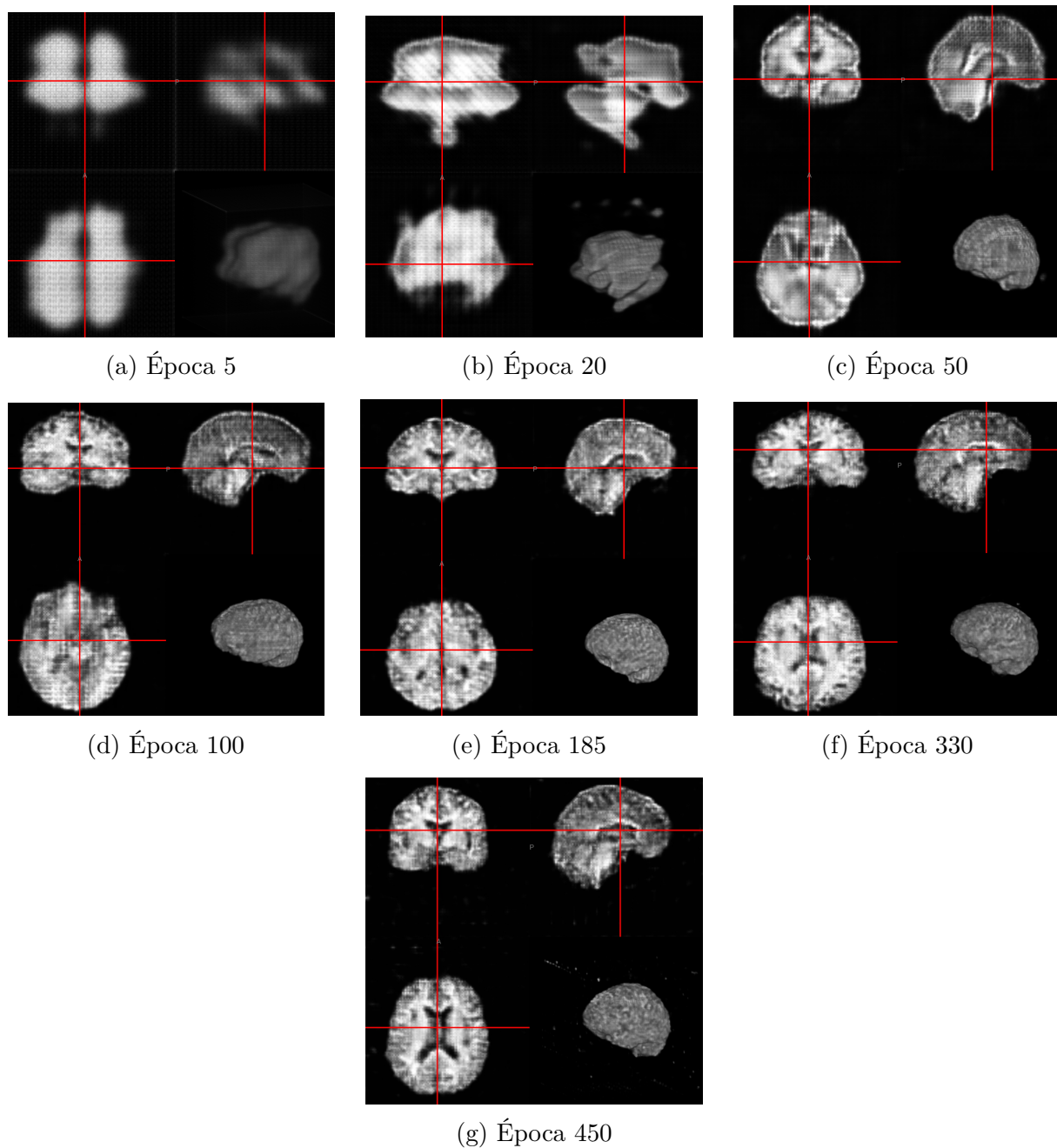


Figura 7.4: Evolución de las muestras sintéticas generadas para la clase CN a lo largo del entrenamiento, mostrando ejemplos en las épocas 5, 20, 50, 100, 185, 330 y 450.

### 7.3.2 Diversidad morfológica de las muestras generadas

Para comprobar que el modelo generativo no se limita a reproducir una única configuración anatómica, sino que es capaz de sintetizar distintos tipos de cerebro dentro de cada clase, se analizan varias realizaciones independientes del generador a partir de vectores latentes distintos. La Figura 7.5 muestra dos ejemplos de volúmenes sintéticos correspondientes a la clase MCI obtenidos tras el entrenamiento del modelo.

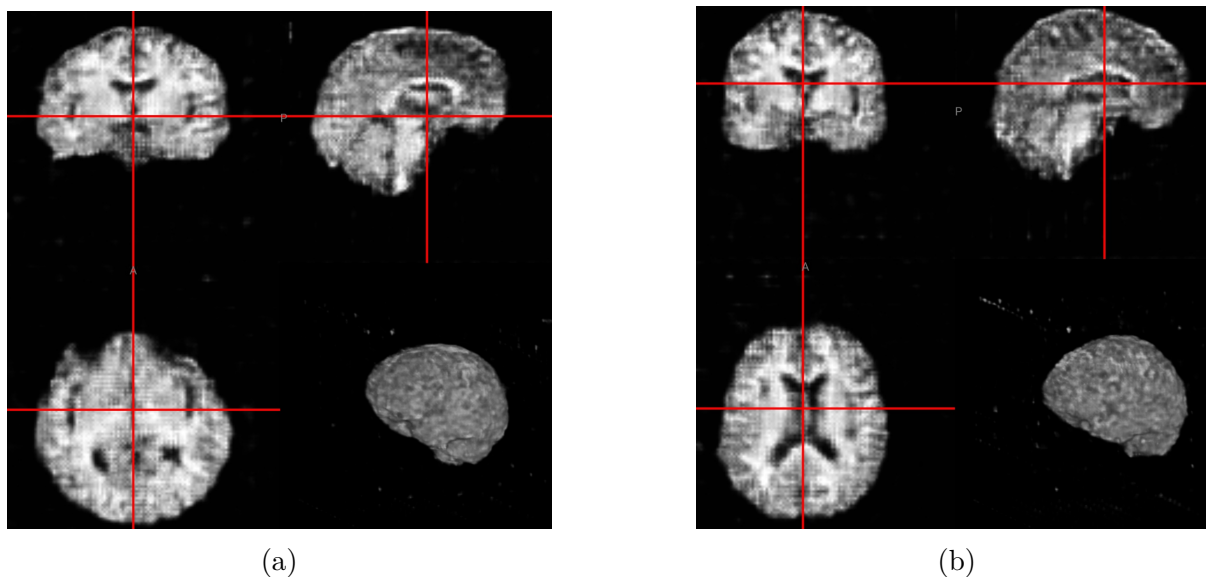


Figura 7.5: Ejemplos de volúmenes sintéticos generados a partir de vectores latentes distintos.

En estos dos ejemplos pueden apreciarse diferencias claras tanto en el tamaño global del encéfalo como en la forma del contorno y en la proporción relativa de las cavidades internas. En el primer caso, el cerebro presenta un volumen ligeramente mayor y una convexidad más marcada de la bóveda craneal, mientras que en el segundo la silueta es más compacta y las estructuras internas ocupan una fracción relativa diferente del espacio intracraneal.

Estas variaciones son coherentes con la diversidad anatómica observada en el dataset real, asociada a factores como el sexo, la complexión corporal o la variabilidad interindividual en el tamaño craneal. El hecho de que el generador produzca volúmenes con diferencias sistemáticas en escala, proporciones y distribución de tejido indica que no está memorizando un único patrón de cerebro, sino que aprende una familia de configuraciones plausibles dentro de esta clase diagnóstica.

Este comportamiento resulta esencial para el objetivo de este trabajo, ya que el propósito del modelo generativo no es replicar exactamente los volúmenes originales, sino proporcionar nuevas muestras anatómicamente coherentes que amplíen la variabilidad del conjunto de entrenamiento. En particular, se plantea incorporar en torno a 30 volúmenes sintéticos por cada grupo (CN, MCI y AD), es decir, aproximadamente un centenar de nuevas imágenes en total, asegurando que dichas muestras no sean prácticamente idénticas entre sí ni copias directas de los sujetos reales. De lo contrario, el uso de datos sintéticos perdería su sentido, al no aportar diversidad efectiva al conjunto de datos ni mejorar la capacidad de generalización del clasificador.

### 7.3.3 Imágenes sintéticas finales

La Figura 7.6 ilustra ejemplos representativos de los volúmenes sintéticos generados para cada una de las clases diagnósticas. En el caso CN (Figura 7.6a) se aprecia un encéfalo con una corteza relativamente continua y una relación compacta entre el parénquima y el espacio ventricular, sin dilataciones llamativas ni adelgazamiento marcado de la superficie cortical.

El ejemplo sintético de MCI (Figura 7.6b) muestra una configuración intermedia: la morfología global se mantiene similar a la de un cerebro cognitivamente normal, pero se observa un ligero aumento del espacio líquido en regiones periventriculares y una mayor irregularidad en la intensidad de la corteza, rasgos compatibles con cambios estructurales incipientes.

Por último, en el volumen generado para la clase AD (Figura 7.6c) se aprecia un incremento más evidente del espacio ventricular y una reducción relativa del tejido cortical, especialmente en regiones temporales y parietales, reflejando un patrón de atrofia más acusado. Aunque las imágenes siguen presentando cierto grado de suavizado y artefactos propios del proceso generativo, estas diferencias cualitativas entre CN, MCI y AD son coherentes con las tendencias morfológicas descritas en la literatura y sugieren que los modelos entrenados han captado, al menos de forma aproximada, las características estructurales que distinguen cada grupo.

Cabe destacar, además, que el modelo correspondiente a la clase CN se ha entrenado con aproximadamente 400 volúmenes reales, mientras que los modelos de MCI y AD disponen de en torno a 200 ejemplos cada uno. Esta diferencia en el tamaño muestral se refleja en que las muestras sintéticas CN presentan, en términos cualitativos, un grado de definición algo superior, pese a que las tres clases alcanzan un nivel de realismo suficiente para los objetivos de este trabajo.

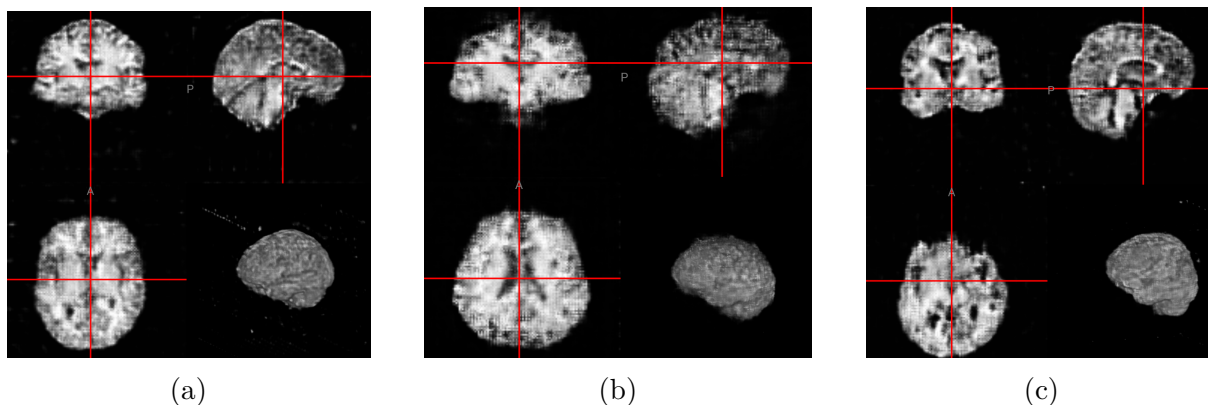


Figura 7.6: Ejemplos representativos de MRIs sintéticos generados para cada clase diagnóstica ( a) CN, b) MCI y c) AD).

### 7.3.4 Denoising de los cerebros sintéticos

Aunque los cerebros generados por los modelos WGAN presentan una anatomía global coherente, muchas de las realizaciones muestran un nivel de ruido apreciable, especialmente visibles en las regiones corticales y en el fondo de la imagen. Antes de incorporar estas imágenes sintéticas al conjunto de entrenamiento del clasificador, se decide aplicar un proceso específico de *denoising* para reducir dicho ruido y homogenizar la textura, de modo que las nuevas muestras se comportasen de forma más similar a las MRI reales preprocesadas.

Inicialmente se evalúan dos enfoques clásicos basados en filtrado gaussiano y *Non-Local Means* (NLM). El primer pipeline aplicaba una máscara de fondo, NLM *slice a slice* con normalización interna y, opcionalmente, un suavizado gaussiano adicional. Aunque este esquema mejoraba parcialmente la apariencia visual de algunas regiones, en otros casos tendía a suavizar en exceso la imagen y a perder contraste fino, especialmente en estructuras pequeñas, por lo que se descartó como solución principal.

El pipeline finalmente adoptado se basa en una combinación de corrección de artefactos de Gibbs, NLM moderado y un realce de nitidez tipo *unsharp masking*. En primer lugar, se aplica la función `gibbs_removal` de DIPY [53] de forma sucesiva a lo largo de los tres ejes del volumen, tras un desplazamiento previo de las intensidades para manejar correctamente los valores negativos. A continuación, se utiliza un NLM más conservador (con un valor de  $h$  proporcional pero no excesivo respecto a la estimación de ruido) calculado plano a plano en el eje axial. Finalmente, se recupera parcialmente el contraste mediante una máscara de realce basada en la diferencia entre la imagen suavizada y la imagen denoised, ajustando el radio ( $\sigma$ ) y la intensidad del realce para evitar la introducción de artefactos.

Sobre esta base se definieron tres configuraciones de parámetros, que dan lugar a tres variantes de volumen denoised (Modelo 1, Modelo 2 y Modelo 3). El Modelo 1 prioriza un aumento del contraste y de la intensidad global, de manera que la imagen resultante parece más limpia pero también algo más artificial, con bordes muy marcados. El Modelo 2 refuerza todavía más este efecto: a simple vista la imagen es muy llamativa, con un contraste elevado y una textura más homogénea, pero algunas estructuras finas empiezan a perderse o a aparecer excesivamente realzadas. El Modelo 3, en cambio, adopta una configuración más conservadora: preserva mejor las intensidades originales, realiza un suavizado menos agresivo y mantiene una textura más cercana a la de las MRI reales, con un equilibrio razonable entre reducción de ruido y fidelidad anatómica.

En términos de parámetros del algoritmo de *denoising*, el Modelo 1 corresponde a una configuración con  $\sigma_{\text{sharp}} = 1,0$  y  $\text{strength} = 1,5$ ; el Modelo 2 emplea  $\sigma_{\text{sharp}} = 1,0$  y

strength = 2,5, lo que acentúa aún más el contraste y el realce de bordes; y el Modelo 3, que es la configuración seleccionada para el resto del trabajo, utiliza valores más conservadores  $\sigma_{\text{sharp}} = 0,8$  y strength = 0,8, priorizando la preservación de la anatomía original frente a un exceso de filtrado.

En la Figura 7.7 se muestra, de izquierda a derecha, la imagen original seguida de las tres variantes denoised. Visualmente, el Modelo 2 puede resultar el más atractivo por su contraste, mientras que el Modelo 3 parece menos espectacular pero más natural. Sin embargo, al contrastar estas primeras impresiones con las métricas de calidad se observa que las soluciones más agresivas no son necesariamente las que mejor respetan la señal original, por lo que la elección final del pipeline se basó en la combinación de la inspección visual y de dichas métricas, y no únicamente en la apariencia subjetiva.

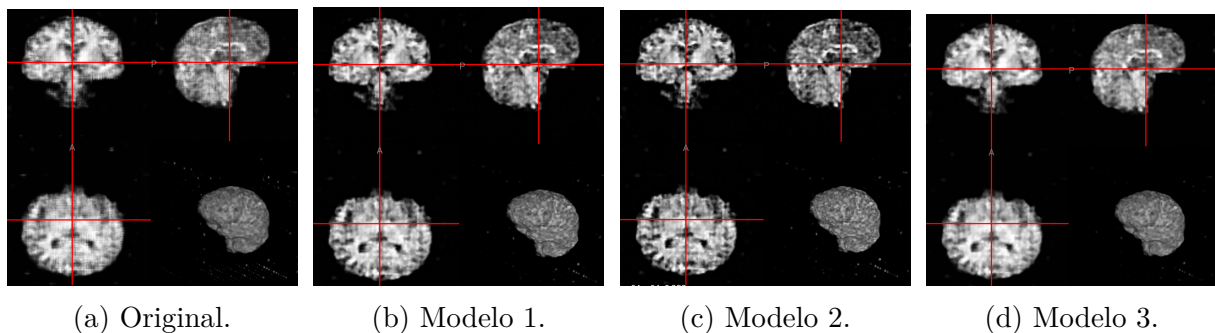


Figura 7.7: Comparación entre la imagen generada original y los tres modelos de denoising evaluados.

Para interpretar los resultados de la Tabla 7.9 se presenta brevemente el significado de cada métrica empleada. El SNR (*Signal-to-Noise Ratio*) cuantifica la relación entre señal y ruido: valores más altos indican imágenes con menos ruido aparente. El PSNR (*Peak Signal-to-Noise Ratio*) compara la imagen original con la denoisada a nivel de píxel y también se interpreta de forma que valores más altos corresponden a una mejor preservación de la señal. El SSIM (*Structural Similarity Index*) evalúa hasta qué punto se mantiene la estructura de la imagen tras el denoising; toma valores entre 0 y 1, siendo 1 el caso ideal de coincidencia estructural (Véase Sección 6.2.1). Por último, el FID (*Fréchet Inception Distance*) mide cuán diferente es la apariencia global de las imágenes denoisadas respecto a las imágenes de referencia; en este caso, cuanto más bajo es el valor, más similares son ambas distribuciones.

La Tabla 7.9 resume los resultados cuantitativos obtenidos para estas tres configuraciones. En todos los casos el SNR original es de 8.2322, de manera que las diferencias se deben exclusivamente al procesamiento. El Modelo 1 reduce el SNR y empeora el FID, mientras que el Modelo 2, pese a su aspecto más contrastado, se aleja todavía más de la distribución de las imágenes reales. El Modelo 3, por el contrario, consigue mejorar lige-

ramente la relación señal–ruido y obtener el mejor FID de las tres variantes, manteniendo al mismo tiempo una apariencia anatómica coherente.

Metricas	Modelo 1	Modelo 2	Modelo 3
SNR original	8.2322	8.2322	8.2322
SNR denoisado	7.3852	6.3957	<b>8.5065</b>
$\Delta$ SNR	-0,8471	-1,8366	<b>+0.2743</b>
PSNR (dB)	30.8977	28.1738	<b>32.7566</b>
SSIM medio	0.9582 $\pm$ 0.0343	0.9477 $\pm$ 0.0366	<b>0.9605 <math>\pm</math> 0.0385</b>
FID	50.0375	58.4406	<b>43.9426</b>

Tabla 7.9: Métricas de evaluación para las tres configuraciones de denoising (Modelos 1–3).

En consecuencia, antes de incorporar los volúmenes sintéticos al conjunto de entrenamiento del modelo de clasificación se aplica un proceso de denoising basado en corrección de Gibbs, NLM moderado y realce de nitidez con los parámetros del Modelo 3. Este pre-procesado reduce el ruido visible, mejora de forma ligera pero consistente las métricas de calidad objetiva y mantiene la apariencia anatómica de las imágenes generadas, lo que resulta especialmente importante al tratarse de un conjunto de datos sintéticos que se utilizará para enriquecer el entrenamiento sin introducir distorsiones sistemáticas en la señal.

# Capítulo 8

## Explicabilidad

### 8.1 Métodos empleados

La explicabilidad en modelos de aprendizaje profundo hace referencia al conjunto de técnicas que permiten entender qué patrones utiliza la red para tomar sus decisiones y cómo contribuye cada región de la entrada a la predicción final [54], [55]. En el contexto de aplicaciones clínicas, esta información es fundamental para evaluar si el modelo se apoya en biomarcadores plausibles desde el punto de vista médico y para generar confianza en sus salidas por parte de los profesionales sanitarios [56].

En este trabajo se ha empleado principalmente Grad-CAM++ (Generalized Gradient-weighted Class Activation Mapping) adaptado a nuestra ResNet34 [57], [58]. Este método extiende Grad-CAM al utilizar una combinación ponderada de las derivadas positivas de los mapas de activación de la última capa convolucional con respecto a la clase de interés, lo que permite obtener mapas de atención más precisos y mejor localizados que el Grad-CAM original, especialmente en escenarios donde pueden coexistir múltiples patrones relevantes en la misma imagen [54], [57], [58]. A partir de los gradientes de forma (1, 512, 4, 4, 4), se obtiene un mapa de atención tridimensional (*Class Activation Map*, CAM) de tamaño  $128 \times 128 \times 128$  que indica, para cada vóxel, su contribución a la probabilidad de la clase predicha. Este CAM se interpola al espacio original del volumen y se superpone sobre la imagen T1 para generar un mapa de calor que resalta las regiones cerebrales más relevantes para la decisión del modelo, siguiendo estrategias similares a las empleadas en trabajos recientes de explicabilidad en Alzheimer [56], [59].

## 8.2 Resultados de explicabilidad

La Figura 8.1 recoge ejemplos de mapas Grad-CAM++ para un sujeto representativo de cada grupo cognitivo (AD, MCI y CN). En todos los casos se selecciona un sujeto del conjunto de test correctamente clasificado por el modelo, de modo que las visualizaciones reflejan las regiones cerebrales que han contribuido de forma más determinante a la predicción correcta. En cada fila de la figura se muestran tres vistas ortogonales (sagital, coronal y axial) del mapa Grad-CAM++ superpuesto sobre la imagen T1 correspondiente. Las zonas en colores cálidos (rojo y amarillo) indican una mayor contribución al aumento de la probabilidad de la clase considerada, mientras que los tonos fríos (azul y verde) señalan regiones con poca o nula influencia en la decisión de la red.

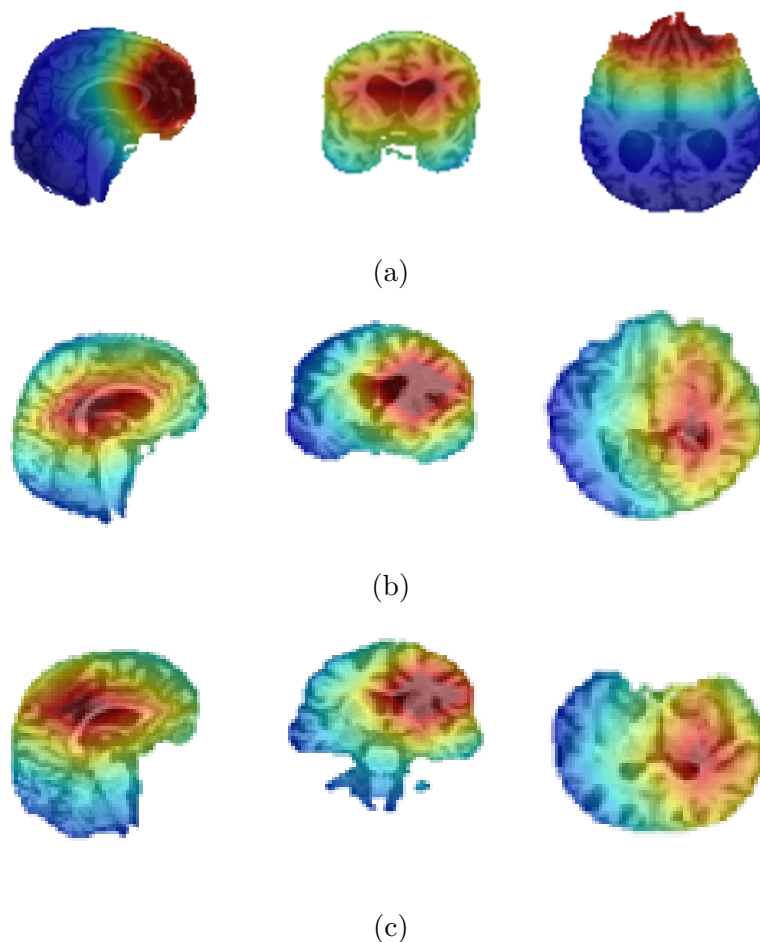


Figura 8.1: Ejemplos de mapas Grad-CAM++ para sujetos del conjunto de test pertenecientes a cada grupo cognitivo: (a) AD, (b) MCI y (c) CN. En cada fila se muestran las vistas sagital, coronal y axial del mapa de atención superpuesto sobre la imagen T1.

En el caso de los sujetos AD (Figura 8.1a), los mapas de activación se distribuyen por gran parte de la corteza y alrededor de los ventrículos laterales, lo que sugiere que

el modelo se apoya en un patrón de atrofia más global coherente con estadios avanzados de la enfermedad. Para los sujetos MCI (Figura 8.1b), la activación se concentra sobre todo en estructuras temporales mediales y regiones adyacentes del lóbulo temporal, en línea con la afectación temprana del hipocampo y la corteza entorrinal descrita en la literatura para el deterioro cognitivo leve. Por último, en los sujetos CN (Figura 8.1c) la activación es más difusa y de menor intensidad, sin focos claros de atrofia, lo que resulta consistente con la ausencia de cambios estructurales marcados esperables en individuos cognitivamente normales.

En conjunto, estos ejemplos sugieren que la ResNet34 está utilizando patrones anatómicos plausibles para distinguir entre CN, MCI y AD: en los casos patológicos el Grad-CAM++ enfatiza regiones clásicamente asociadas a neurodegeneración, mientras que en los sujetos sanos no aparecen focos de activación espurios en estructuras irrelevantes (bordes del cráneo, ruido fuera del cerebro).

### 8.2.1 Explicabilidad en modelos con y sin datos sintéticos

Además del análisis de explicabilidad sobre un único modelo, se ha estudiado cómo cambia el patrón de atención cuando se aplica Grad-CAM++ a dos variantes del clasificador: (i) una ResNet34 entrenada únicamente con imágenes reales y (ii) la misma arquitectura entrenada con una combinación de imágenes reales y volúmenes sintéticos generados por el WGAN-3D. El objetivo es comprobar si la inclusión de datos sintéticos modifica de manera apreciable las regiones que cada modelo considera más relevantes para tomar sus decisiones.

La Figura 8.2 muestra un ejemplo de este experimento para un sujeto cognitivamente normal (CN). En la fila (a) se representan las vistas sagital, coronal y axial correspondientes al modelo entrenado sólo con datos reales, mientras que en la fila (b) se muestran las mismas vistas para el modelo entrenado con datos reales y sintéticos. En ambos casos se ha utilizado la misma imagen T1 y se han seleccionado las slices centrales, de forma que la comparación entre patrones de activación sea directa. Para este sujeto, ambos modelos predicen correctamente la clase CN, con probabilidades cercanas a 1 para la clase correcta y valores muy bajos para MCI y AD.

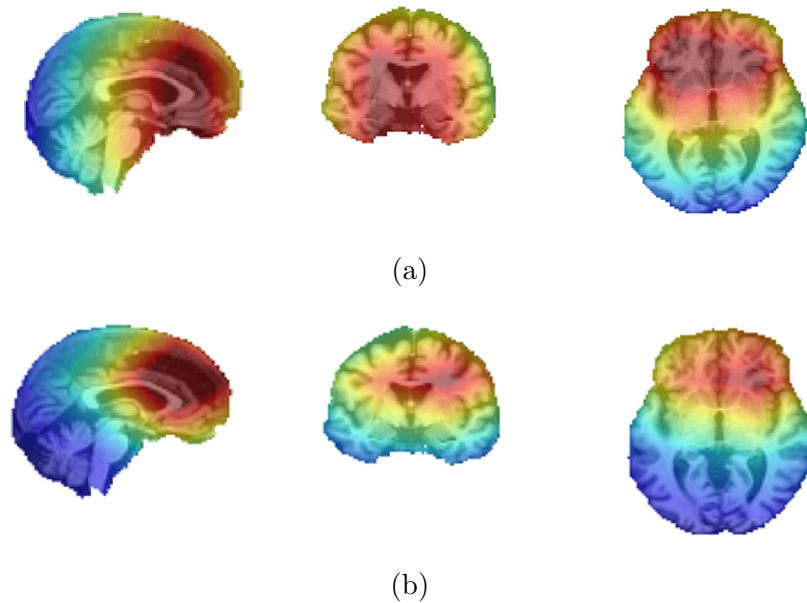


Figura 8.2: Mapas Grad-CAM++ para un sujeto cognitivamente normal (CN) obtenidos con dos modelos distintos. (a) Modelo entrenado sólo con imágenes reales. (b) Modelo entrenado con imágenes reales y volúmenes sintéticos. En ambos casos se muestran las proyecciones sagital, coronal y axial con el mapa de activación superpuesto sobre la misma imagen T1.

En el caso CN se observa que los dos modelos centran la mayor parte de la activación en zonas relativamente similares, con un patrón difuso y de baja intensidad repartido por la corteza. Esto es coherente con el hecho de que se trata de un sujeto sin alteraciones marcadas: el modelo no necesita apoyarse en cambios estructurales muy localizados para confirmar que la imagen pertenece al grupo normal, y la incorporación de datos sintéticos no introduce focos de activación claramente nuevos o extraños.

En el caso de Alzheimer (AD) (Figura 8.3), se aprecia que los dos modelos generan mapas de activación distintos a pesar de partir de la misma imagen T1. En la Figura 8.3a, correspondiente al modelo entrenado sólo con datos reales, el Grad-CAM++ muestra un patrón de activación bastante extendido y relativamente uniforme por gran parte de la corteza, con colores cálidos ocupando buena parte de las tres vistas. Esto sugiere que este modelo reparte la atención por muchas regiones del cerebro a la vez para justificar la predicción de AD.

En la Figura 8.2b, que representa el modelo entrenado con datos reales y sintéticos, el mapa de calor aparece más concentrado. Las zonas de alta activación se desplazan hacia regiones más concretas y la intensidad disminuye en áreas donde antes el mapa era casi completamente rojo. Dicho de otro modo, el segundo modelo parece focalizarse en un conjunto más reducido de zonas relevantes, mientras que el primero distribuye su atención de forma más difusa.

Esta diferencia entre un patrón más “global” (modelo sólo real) y otro más “focal” (modelo real+syntético) indica que la inclusión de datos sintéticos no cambia por completo las regiones que el modelo considera importantes, pero sí influye en cómo reparte la atención: el modelo aumentado tiende a concentrarse en menos áreas con activación alta y a reducir el peso de regiones que podrían ser menos informativas para la decisión.

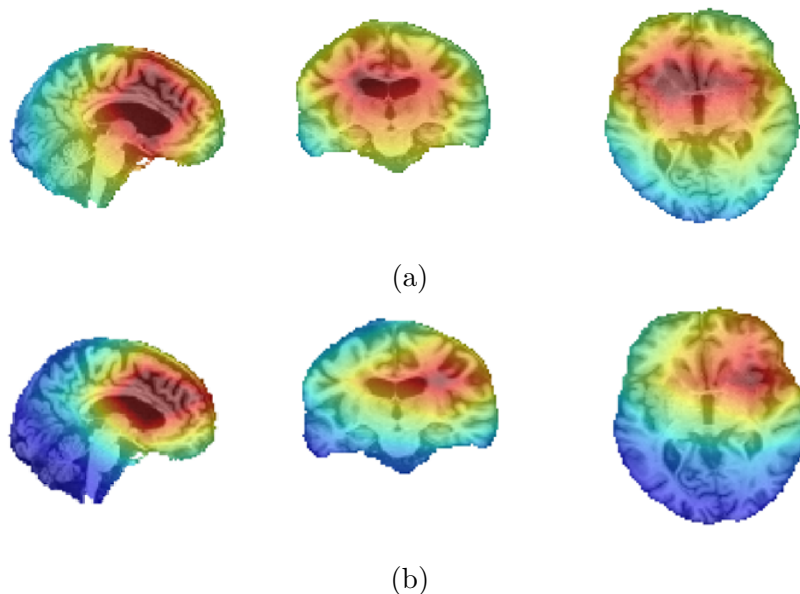


Figura 8.3: Mapas Grad-CAM++ para un sujeto con Alzheimer (AD) obtenidos con dos modelos distintos. (a) Modelo entrenado sólo con imágenes reales. (b) Modelo entrenado con imágenes reales y volúmenes sintéticos. Se representan las proyecciones sagital, coronal y axial con el mapa de activación superpuesto sobre la misma imagen T1.

En conjunto, estas comparaciones sugieren que añadir volúmenes sintéticos al entrenamiento no cambia de forma radical las zonas del cerebro en las que el modelo se fija para distinguir entre CN y AD: las regiones destacadas son, en esencia, las mismas. Grad-CAM++ produce patrones de activación similares para ambos modelos cuando se evalúan sobre la misma imagen, aunque el modelo entrenado con datos sintéticos tiende a concentrar más la atención en determinadas áreas. Esto indica que el aumento con datos sintéticos no introduce patrones de atención inesperados, sino que refuerza y afina el comportamiento que ya se observa con un entrenamiento basado únicamente en datos reales, probablemente porque este segundo modelo ha visto un mayor número de ejemplos durante el entrenamiento y ha aprendido a focalizar mejor las regiones más informativas para la clasificación.

## 8.2.2 Comparación de los mapas Grad-CAM++ en modelos multiclase y binarios

En esta sección se analizan los mapas Grad-CAM++ generados para un mismo sujeto del grupo AD al aplicar el modelo multiclase (AD/MCI/CN) y el modelo binario (AD vs. CN). Grad-CAM++ permite visualizar qué regiones del cerebro contribuyen en mayor medida a la predicción de la red, proporcionando mapas de activación específicos de clase a partir de los gradientes en las últimas capas convolucionales.[60]

En la Figura 8.4a se muestran las proyecciones sagital, coronal y axial del mapa de activación correspondiente al modelo multiclase. En este caso, las zonas de mayor intensidad se distribuyen de forma relativamente amplia por la corteza, con especial énfasis en regiones temporales y parietales, coherentes con la atrofia característica de la enfermedad de Alzheimer descrita en la literatura. Esta distribución algo más difusa es consistente con el hecho de que el modelo debe separar simultáneamente tres categorías (AD, MCI y CN), por lo que su atención se reparte entre patrones más heterogéneos asociados a los diferentes estadios de la enfermedad.

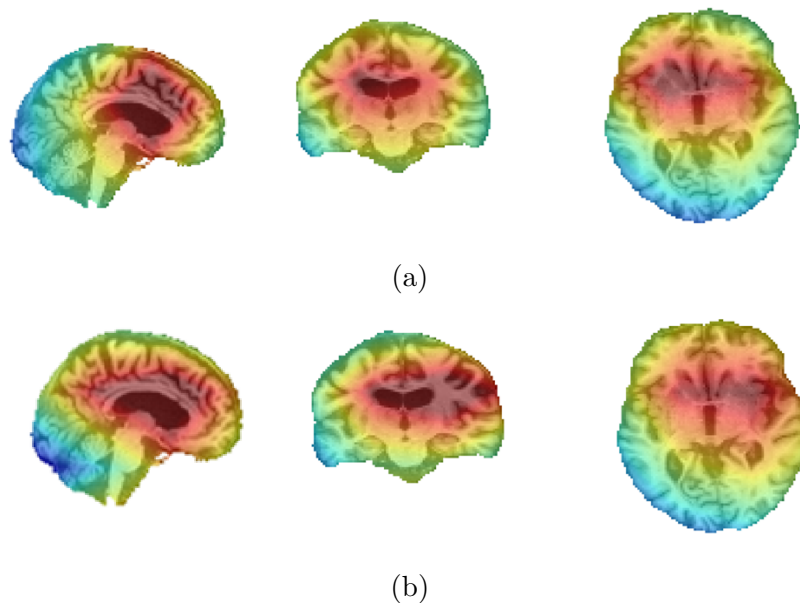


Figura 8.4: Mapas Grad-CAM++ correspondientes a un mismo sujeto con Alzheimer (AD) obtenidos para dos configuraciones de entrenamiento distintas. (a) Modelo multiclase (AD/MCI/CN) entrenado exclusivamente con volúmenes reales. (b) Modelo binario (AD vs. CN) entrenado con la combinación de volúmenes reales y sintéticos. Se muestran las proyecciones sagital, coronal y axial con el mapa de activación superpuesto sobre la imagen T1.

Por su parte, la Figura 8.4b recoge los mapas Grad-CAM++ obtenidos con el modelo binario entrenado para distinguir únicamente entre AD y CN. Visualmente, las regiones

resaltadas son muy similares a las del modelo multiclase, lo que indica que ambos modelos se apoyan en patrones anatómicos comparables. No obstante, en el modelo binario las activaciones tienden a concentrarse ligeramente más en las zonas más discriminativas entre ambos grupos, con un realce algo más compacto en estructuras medial-temporales y regiones corticales donde se esperan cambios estructurales relevantes en AD.[60], [61]

Al tratarse de una tarea de clasificación más simple, el modelo binario puede especializarse mejor en los patrones que diferencian de forma directa a AD de sujetos sanos, lo que se traduce en mapas de atención ligeramente más focalizados y coherentes anatómicamente. En conjunto, la comparación sugiere que la formulación binaria no modifica de forma sustancial las regiones cerebrales en las que la red se fija para detectar AD, pero sí favorece una atención algo más concentrada, mientras que la formulación multiclase tiende a mostrar activaciones más extendidas. Estas diferencias son coherentes con lo observado en otros trabajos de clasificación de Alzheimer a partir de neuroimagen, donde los modelos binarios suelen mostrar focos de activación más localizados y los modelos multiclase reflejan la mayor complejidad y solapamiento entre estadios intermedios como MCI.[61]

### 8.3 Implicaciones clínicas

Desde una perspectiva clínica, disponer de mapas de importancia como los generados por Grad-CAM++ aporta transparencia a las decisiones del modelo y reduce la sensación de “caja negra”. Los neurorradiólogos pueden ver qué regiones han sustentado la predicción (por ejemplo, MCI o AD) y comprobar si coinciden con los hallazgos que ellos mismos identificarían en la resonancia, lo que facilita la detección de errores sistemáticos y favorece la aceptación del sistema como herramienta de apoyo. Además, el hecho de que estos mapas muestren patrones anatómicamente plausibles tanto en imágenes reales como en volúmenes sintéticos indica que los datos generados han capturado rasgos estructurales relevantes de la enfermedad, en lugar de introducir artefactos, reforzando así la confianza en el uso de datos sintéticos para entrenar modelos clínicos.

# Capítulo 9

## Discusión

### 9.1 Interpretación de resultados

En el ámbito de la imagen médica, el desarrollo de modelos de *aprendizaje profundo* realmente útiles en la práctica clínica se ve limitado por varios factores estructurales: el tamaño reducido de muchos conjuntos de datos, el desequilibrio entre clases (muchos controles sanos frente a relativamente pocos pacientes), las restricciones de privacidad que dificultan compartir información entre centros y el coste elevado de obtener anotaciones expertas [62]. Como consecuencia, entrenar redes profundas con capacidad de generalización resulta complejo, y el progreso en herramientas automáticas para diagnóstico o estratificación de pacientes es más lento de lo deseable. En el caso concreto de la enfermedad del Alzheimer, estos problemas se acentúan al tratar de distinguir fases intermedias como el deterioro cognitivo leve (MCI), donde las diferencias anatómicas son sutiles y la cantidad de datos disponibles suele ser menor que para los grupos claramente patológicos o claramente sanos [63].

En este contexto, la estrategia seguida en este trabajo combina dos líneas complementarias para mitigar la escasez y el desbalanceo de datos: por un lado, el entrenamiento de un clasificador 3D basado en ResNet34, ajustado específicamente al tamaño y características del conjunto de MRI estructurales disponible; por otro, la generación de volúmenes T1 sintéticos mediante modelos WGAN-3D, con el objetivo de emplearlos como *data augmentation* controlado. A partir del conjunto original de 882 volúmenes (458 CN, 231 MCI, 193 AD), se entrenó en primer lugar una ResNet34 únicamente con imágenes reales, que sirve como modelo de referencia para evaluar el impacto de los datos sintéticos.

Para incorporar las nuevas imágenes generadas, se construyó un fichero de etiquetas ampliado, `labels_3classes_with_synthetic`, a partir del CSV original `labels_3classes`.

En este nuevo fichero se añadieron 90 volúmenes sintéticos adicionales, 30 por cada grupo diagnóstico (CN, MCI y AD), asignándoles la misma etiqueta de clase que a los sujetos reales correspondientes y dejando marcados con valores neutros (por ejemplo, “X”) los campos sensibles como sexo y edad. De este modo, el dataset de entrenamiento se amplió de forma explícita y trazable, manteniendo una estructura de metadatos homogénea entre imágenes reales y sintéticas. En la tabla 9.1 se muestra un ejemplo del dataset final.

Subject	Group	Sex	Age	Acq Date
136_S_1227	MCI	F	65	02/21/2007
136_S_0579	AD	F	66	07/10/2006
brain_gen_1	CN	X	X	06/03/2026
brain_gen_35	AD	X	X	06/07/2026
brain_gen_87	MCI	X	X	06/11/2026

Tabla 9.1: Ejemplo del dataset final de sujetos combinando las imágenes originales de ADNI con las imágenes sintéticas generadas por la WGAN.

La Figura 9.1b muestra la evolución de la pérdida y la *exactitud* durante el entrenamiento de la ResNet34 con el conjunto mixto de datos reales y sintéticos. La *pérdida de entrenamiento* desciende de manera sostenida y tiende a estabilizarse en las últimas épocas, mientras que la *pérdida de validación* presenta oscilaciones moderadas pero mantiene una tendencia global decreciente, sin indicios claros de divergencia. De forma coherente, la *exactitud de entrenamiento* aumenta progresivamente hasta acercarse al 100 %, y la *exactitud de validación* alcanza valores en torno al 75–77 %, situándose de forma estable por encima del 70 % a partir de la segunda mitad del entrenamiento. Este comportamiento sugiere que el modelo es capaz de aprovechar el aumento de datos sin caer en un sobreajuste extremo, manteniendo una brecha controlada entre las curvas de entrenamiento y validación. A partir de estas curvas se seleccionó como modelo final la época correspondiente al mejor compromiso entre *pérdida de validación* y *exactitud de validación*, que coincide con una *Val Acc* máxima cercana al 76–77 %.

La Tabla de métricas 9.2 y la matriz de confusión 9.3 muestran el rendimiento de la ResNet34 entrenada con el conjunto mixto de datos reales y sintéticos. En el conjunto de test (146 volúmenes, distribuidos en 34 AD, 73 CN y 39 MCI), el modelo alcanza valores de precisión y *recall* elevados para la clase CN (0.88 y 0.86, respectivamente), lo que se traduce en un F1 de 0.87 y en un número muy reducido de falsos positivos y falsos negativos para este grupo. En las clases patológicas, el comportamiento es algo más modesto pero razonablemente equilibrado: para AD se obtiene una precisión de 0.64 y una sensibilidad de 0.74 (F1 = 0.68), mientras que para MCI la precisión es de 0.74 y la sensibilidad de 0.67 (F1 = 0.70), indicando que el modelo es capaz de detectar una proporción sustancial de sujetos en fases intermedias sin disparar en exceso los falsos positivos.

Clase	Precisión	<i>Recall</i>	F1
AD	0.64	0.74	0.68
CN	0.88	0.86	0.87
MCI	0.74	0.67	0.70

Tabla 9.2: Métricas de test por clase para la ResNet34 con datos reales y sintéticos.

<b>Real / Pred</b>	AD	CN	MCI
AD	25	5	4
CN	5	63	5
MCI	9	4	26

Tabla 9.3: Matriz de confusión en test para la ResNet34 con datos reales y sintéticos.

La matriz de confusión de la Tabla 9.3 permite analizar con más detalle cómo se distribuyen los errores entre clases. En la diagonal se observa que el modelo acierta 25 de los 34 casos de AD, 63 de los 73 CN y 26 de los 39 MCI, lo que confirma el buen rendimiento global ya reflejado en las métricas de F1. Los errores más frecuentes aparecen entre las clases patológicas: 9 sujetos MCI son etiquetados como AD y 4 AD se confunden con MCI, lo que sugiere que el modelo tiende a desplazar parte de los casos intermedios hacia la categoría más grave, en línea con la continuidad clínica entre MCI y AD. En cambio, las confusiones entre CN y las clases patológicas son menos numerosas (5 CN clasificados como AD y 5 como MCI, y 4 MCI clasificados como CN), lo que indica que la red distingue razonablemente bien entre cerebros sanos y patológicos, concentrando la incertidumbre principalmente en la frontera entre MCI y AD.

### 9.1.1 Comparación de modelos con y sin datos sintéticos

En las secciones anteriores se han presentado por separado los resultados del clasificador ResNet34 entrenado únicamente con volúmenes T1 reales y del mismo modelo tras incorporar volúmenes sintéticos generados mediante WGAN-3D. En esta subsección se comparan de forma directa ambas configuraciones, analizando primero las curvas de entrenamiento y validación, después las métricas globales y por clase en el conjunto de test y, finalmente, las matrices de confusión. El objetivo es evaluar hasta qué punto el aumento de datos sintéticos contribuye a mejorar la capacidad de generalización del clasificador y, en particular, su rendimiento en la detección de sujetos con deterioro cognitivo leve (MCI).

Las curvas de la Figura 9.1 comparan directamente el comportamiento de la ResNet34 entrenada solo con datos reales (panel 9.1a) y de la misma arquitectura cuando se incorporan volúmenes sintéticos (panel 9.1b). En ambos casos la *pérdida de entrenamiento* desciende de forma monótona y la *exactitud de entrenamiento* se aproxima al 100 %, pero las trayectorias de validación muestran diferencias importantes en términos de estabilidad y sobreajuste.

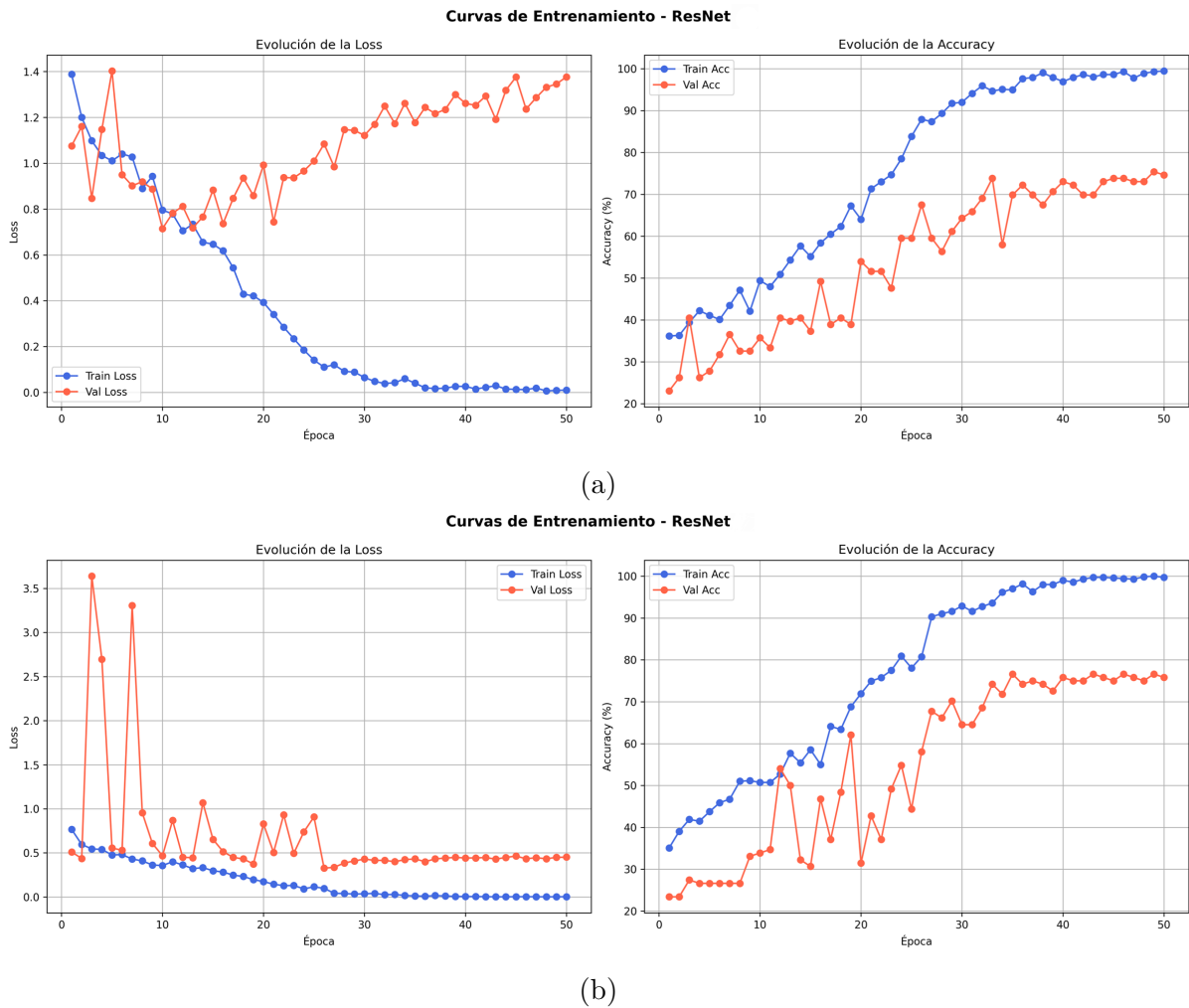


Figura 9.1: Curvas de pérdida (izquierda) y *exactitud* (derecha) en entrenamiento y validación para la ResNet34, comparando el modelo entrenado solo con datos reales (a) y el modelo entrenado con datos reales más volúmenes sintéticos (b), ambos modelos entrenados con los mismos parámetros y la misma estructura de ResNet34.

En el modelo entrenado únicamente con imágenes reales (Figura 9.1a), la *pérdida de entrenamiento* disminuye de forma clara y sostenida, mientras que la *pérdida de validación* desciende en las primeras épocas pero pronto se estabiliza e incluso muestra una tendencia ligeramente creciente, con oscilaciones relativamente amplias. De manera coherente, la *exactitud de entrenamiento* se aproxima rápidamente a valores superiores al 95 %, mientras que la *exactitud de validación* se mantiene en torno al 70–74 %, con una brecha cada vez mayor entre ambas curvas. Este patrón es característico de un cierto grado de sobreajuste: el modelo acaba representando muy bien el conjunto de entrenamiento, pero su capacidad de generalización a datos no vistos es más limitada.

Cuando se añaden los volúmenes sintéticos (Figura 9.1b), la dinámica de entrenamiento cambia de forma apreciable. La *pérdida de entrenamiento* sigue disminuyendo hasta valores cercanos a cero, pero la *pérdida de validación* desciende más de manera más prolon-

gada y termina estabilizándose en niveles inferiores, con oscilaciones menos pronunciadas que en el caso sin aumento sintético. Del mismo modo, la *exactitud de validación* se sitúa de forma consistente por encima del 70 % y alcanza picos alrededor del 76–77 %, con una separación algo más contenida respecto a la *exactitud de entrenamiento*. En conjunto, estas curvas sugieren que el aumento con datos generados aporta variabilidad adicional que ayuda a regularizar el aprendizaje: aunque el modelo sigue siendo capaz de memorizar el conjunto de entrenamiento, la mejor evolución de la pérdida y de la *exactitud* en validación indica un sobreajuste menos acusado y una capacidad de generalización superior.

Desde el punto de vista cuantitativo, la ResNet34 entrenada solo con imágenes reales alcanza en el conjunto de test una *exactitud* del 73.68 %, como se aprecia en la Tabla 9.4. Al incorporar las 90 imágenes sintéticas generadas por los WGAN-3D, la misma arquitectura alcanza una *exactitud de test* del 78.08 % (Tabla 9.5). Es decir, se obtiene una mejora de aproximadamente 4–5 puntos porcentuales en *exactitud*, acompañada de un incremento paralelo en las métricas basadas en F1, lo que indica que el aumento sintético no solo ayuda al modelo a acertar más, sino también a mantener un equilibrio más favorable entre precisión y sensibilidad.

La comparación detallada por clase revela cambios importantes en el comportamiento del clasificador. En el modelo entrenado únicamente con datos reales, las clases AD y CN muestran métricas muy equilibradas, con F1 de 0.77 y 0.78 respectivamente, gracias a combinaciones de precisión y *recall* relativamente altas (Tabla 9.4). Sin embargo, la clase MCI se ve claramente penalizada: aunque la precisión es elevada (0.80), la sensibilidad se reduce hasta 0.46, lo que se traduce en un F1 de solo 0.58. En la práctica, esto significa que el modelo es muy conservador a la hora de etiquetar MCI: cuando predice esta clase suele acertar, pero deja sin detectar una proporción considerable de sujetos con deterioro cognitivo leve, que acaban siendo asignados a CN o AD.

Clase	Precision	Recall	F1
CN	0.72	<b>0.84</b>	<b>0.78</b>
MCI	<b>0.80</b>	0.46	0.58
AD	0.73	0.83	0.77

Tabla 9.4: Métricas de test para el modelo con datos reales

Clase	Precisión	<i>Recall</i>	F1
CN	<b>0.88</b>	<b>0.86</b>	<b>0.87</b>
MCI	0.74	0.67	0.70
AD	0.64	0.74	0.68

Tabla 9.5: Métricas de test por clase para el modelo con datos reales y sintéticos.

Cuando se incorporan las imágenes sintéticas, el balance entre clases cambia de forma notable (Tabla 9.5). La clase CN es la que más mejora: su precisión aumenta de 0.72 a 0.88 y el F1 pasa de 0.78 a 0.87, manteniendo una sensibilidad muy similar (de 0.84 a 0.86), lo que indica que el modelo comete menos falsos positivos para controles sanos sin sacrificar

la capacidad de detectarlos. La clase AD sufre una ligera degradación, bajando de un F1 de 0.77 a 0.68 debido sobre todo a la disminución de la precisión (de 0.73 a 0.64), aunque conserva un *recall* relativamente alto (0.83 frente a 0.74), de modo que sigue identificando la mayoría de los casos de Alzheimer. La mejora más relevante se observa de nuevo en MCI: el F1 se incrementa de 0.58 a 0.70 gracias a un aumento claro de la sensibilidad (de 0.46 a 0.67), incluso aceptando una ligera reducción de la precisión (de 0.80 a 0.74). En términos clínicos, esto implica que el modelo con datos sintéticos identifica muchos más pacientes con deterioro cognitivo leve, aunque a costa de etiquetar como MCI algunos casos que no lo son.

En conjunto, estos resultados muestran que las imágenes sintéticas no solo elevan la *precisión* global, sino que corrigen en gran medida el desequilibrio inicial entre clases. El clasificador deja de infradetectar la etapa intermedia de la enfermedad (MCI) y mejora su rendimiento en CN, aceptando un pequeño sacrificio en la precisión de AD. Dado que el objetivo del sistema es priorizar la detección temprana y reducir al mínimo los falsos negativos en MCI, este nuevo equilibrio entre precisión y *recall* puede considerarse clínicamente más deseable que el obtenido con el modelo entrenado exclusivamente con datos reales.

Las matrices de confusión de la Tabla 9.6 y la Tabla 9.7 refuerzan esta lectura y permiten cuantificar mejor el impacto del aumento sintético sobre la clase MCI, que es la más difícil de predecir en este problema. Sin imágenes sintéticas, la ResNet34 acierta 24 de 29 casos de AD, 58 de 69 CN y únicamente 16 de 35 MCI; además, comete 17 errores en los que sujetos MCI se clasifican como CN y solo 2 MCI se confunden con AD (Tabla 9.6). Este patrón revela que el modelo tiende a predecir una gran parte de los pacientes con deterioro cognitivo leve hacia el grupo de controles, infraestimando la presencia de MCI y dejando sin detectar casi la mitad de los casos reales.

Cuando se incorporan los volúmenes sintéticos, la distribución de aciertos y errores cambia de forma notable (Tabla 9.7). El modelo pasa a acertar 25 de 34 AD, 63 de 73 CN y, sobre todo, 26 de 39 MCI, lo que supone un incremento de 10 aciertos adicionales en esta última clase. Al mismo tiempo, las confusiones MCI→CN se reducen drásticamente de 17 a solo 4 casos, mientras que aumentan las transiciones MCI→AD (de 2 a 9), desplazando la incertidumbre desde la frontera MCI–CN hacia la frontera MCI–AD, más coherente con la progresión clínica de la enfermedad. En términos de precisión para MCI, esto implica que el clasificador comete menos falsos negativos y se muestra mucho más capaz de identificar correctamente a los sujetos con deterioro cognitivo leve, que constituyen el objetivo más crítico del estudio.

Real/Pred	CN	MCI	AD
CN	<b>58</b>	4	7
MCI	<b>17</b>	16	2
AD	5	0	<b>24</b>

Tabla 9.6: Matriz de confusión en test para el modelo con datos reales.

Real / Pred	CN	MCI	AD
CN	<b>63</b>	5	5
MCI	4	<b>26</b>	9
AD	5	4	<b>25</b>

Tabla 9.7: Matriz de confusión en test para el modelo con datos reales y sintéticos.

En conjunto, la comparación indica que el aumento sintético introduce un compromiso distinto: el modelo pierde algo de precisión en AD, pero mejora de manera clara capacidad para distinguir MCI de CN y mantener un alto rendimiento en CN. Desde la perspectiva del cribado y del diagnóstico temprano, este cambio es especialmente relevante, ya que la correcta identificación de MCI es el escenario más desafiante y al mismo tiempo el más valioso desde el punto de vista clínico. Así, el uso de imágenes sintéticas no solo incrementa la exactitud global en test, sino que corrige en gran medida la tendencia inicial a infradetectar MCI y mejora la precisión y la sensibilidad asociadas a esta clase, acercando el comportamiento del clasificador a lo que sería deseable en un entorno real de apoyo a la toma de decisiones.

## 9.2 Comparación con el estado del arte

En esta sección se comparan críticamente los resultados obtenidos con la ResNet34 entrenada con datos reales y sintéticos con los trabajos revisados en el capítulo de estado del arte que abordan la clasificación conjunta de sujetos AD, MCI y CN a partir de MRI estructural. La atención se centra especialmente en el rendimiento para la clase MCI, que constituye el principal cuello de botella debido a la superposición de patrones con CN y AD y a la escasez de muestras etiquetadas [64], [65].

Diversos estudios recientes que utilizan redes 3D convolucionales y descriptores radiómicos sobre T1-weighted MRI obtienen *exactitudes* globales entre el 75 % y el 85 % para la clasificación multiclase AD–MCI–CN, con un rendimiento sensiblemente inferior en la clase MCI, cuyo F1 suele situarse por debajo de 0.70 [64], [65]. En muchos casos, la matriz de confusión revela precisamente el mismo patrón observado en nuestro modelo sin aumento sintético: una tendencia a clasificar una fracción importante de los sujetos MCI como CN, lo que se traduce en una sensibilidad reducida para esta categoría intermedia [64]. Frente a este contexto, la *exactitud de test* del 78.08 % obtenida en este trabajo con la ResNet34 aumentada y, sobre todo, el F1 de 0.70 alcanzado para MCI, resultan competitivos teniendo en cuenta el tamaño limitado del conjunto de entrenamiento y el

uso exclusivo de MRI estructural como modalidad de entrada [64], [65].

Un aspecto especialmente relevante es que esta precisión en MCI se logra a pesar de la escasez de datos, situación que suele conducir a un sesgo marcado hacia CN y AD [65], [66]. En nuestro caso, la incorporación de solo 90 volúmenes sintéticos produce una mejora aproximada de 5% en *exactitud* respecto al mismo modelo entrenado únicamente con datos reales, y se traduce en un incremento sustancial del F1 de MCI (de 0.58 a 0.70) y en una reducción muy notable de las confusiones MCI→CN en la matriz de confusión. Esta combinación de resultados sugiere que el aumento sintético no solo compensa la falta de muestras reales, sino que guía al clasificador hacia una representación más discriminativa de la etapa de deterioro cognitivo leve [66], [67], [68].

Hasta donde alcanza nuestra revisión, no se han encontrado trabajos que evalúen de forma explícita una hipótesis idéntica a la planteada en este estudio: demostrar que un número relativamente pequeño de volúmenes sintéticos 3D, generados de manera específica para un conjunto multiclase AD–MCI–CN, es capaz de mejorar de forma significativa el rendimiento del modelo precisamente en la clase MCI, manteniendo al mismo tiempo una *exactitud* global competitiva [66], [69]. La mayoría de propuestas basadas en GAN se centran en aumentar la robustez general del clasificador o en mejorar el rendimiento global, pero rara vez analizan de manera detallada el efecto del aumento sobre la clase MCI dentro de un escenario multiclase completo [67], [68]. En este sentido, los resultados presentados aquí aportan evidencia empírica novedosa a favor de dicha hipótesis: con tan solo 90 imágenes sintéticas adicionales se consigue un incremento de alrededor del 5% en la *exactitud* de test y una mejora clara en las métricas asociadas a MCI, lo que respalda de forma sólida la utilidad de este tipo de aumento de datos en contextos donde el número de estudios disponibles es limitado [66], [69].

# Capítulo 10

## Conclusiones y trabajo futuro

### 10.1 Conclusiones

En este trabajo se desarrolla un flujo completo para la clasificación automática de sujetos con enfermedad de Alzheimer (AD), deterioro cognitivo leve (MCI) y controles sanos (CN) a partir de volúmenes T1 de resonancia magnética estructural. El flujo de trabajo propuesto abarca desde el preprocesamiento de las imágenes y la generación de datos sintéticos mediante modelos WGAN-3D hasta el entrenamiento de un clasificador profundo basado en ResNet34 y el análisis de explicabilidad mediante mapas de activación. En conjunto, los resultados obtenidos responden de forma positiva a los objetivos planteados y demuestran que es posible alcanzar un rendimiento competitivo en un escenario multiclase AD–MCI–CN a pesar del tamaño limitado del conjunto de datos disponible.

En cuanto al preprocesamiento, se implementa una cadena estandarizada que incluye registro espacial, normalización de intensidad, recorte y redimensionado de los volúmenes, así como estrategias de equilibrado de las clases. Este flujo de trabajo homogeneiza los datos procedentes de diferentes sujetos y garantiza que el modelo profundo reciba entradas coherentes, reduciendo la variabilidad no relevante y facilitando la extracción de patrones estructurales asociados a cada estadio de la enfermedad. La calidad de este preprocesamiento se refleja en la estabilidad de las curvas de entrenamiento y en la ausencia de artefactos evidentes en los volúmenes generados y clasificados.

Respecto al modelo de clasificación, la arquitectura ResNet34 muestra un compromiso adecuado entre capacidad representacional y riesgo de sobreajuste en el contexto de un conjunto de datos moderado. Entrenada únicamente con imágenes reales, la red alcanza una *exactitud* en test del 73.68 %, con valores de F1 macro y ponderado en torno a 0.71–0.73, lo que confirma su capacidad para discriminar entre AD, MCI y CN y pone

de manifiesto, al mismo tiempo, las dificultades específicas en la predicción de MCI. La matriz de confusión del modelo base evidencia que una fracción sustancial de sujetos con MCI se clasifica erróneamente como CN, reproduciendo el patrón descrito en la literatura, donde la clase intermedia suele ser la más problemática.

Un segundo bloque de contribuciones se centra en la generación de imágenes sintéticas mediante WGAN-3D. Se entrenan redes generativas adversarias capaces de producir volúmenes sintéticos realistas para cada una de las clases, se evalúa cualitativamente la similitud visual con las imágenes reales y se comprueba que no se introducen artefactos clínicamente inverosímiles. Estos volúmenes sintéticos se utilizan posteriormente como aumento de datos específico para la tarea de clasificación, con el objetivo de mitigar la escasez de muestras y el desbalance de clases característicos de los conjuntos clínicos de neuroimagen. De este modo, el trabajo no solo entrena un modelo discriminativo, sino que explora también la vertiente generativa como herramienta para mejorar el rendimiento en escenarios de datos limitados.

La explicabilidad del clasificador se aborda mediante técnicas basadas en mapas de activación, como Grad-CAM, que permiten visualizar las regiones cerebrales que contribuyen de forma más relevante a las decisiones del modelo. El análisis de estos mapas muestra activaciones concentradas en estructuras compatibles con los biomarcadores descritos para AD y MCI (por ejemplo, regiones temporales mediales y áreas corticales de asociación), mientras que en los controles sanos la activación es más difusa o de menor intensidad. Estas observaciones indican que la red aprende patrones anatómicamente plausibles y refuerzan la confianza en que las decisiones del modelo no se basan en artefactos espurios, en línea con las recomendaciones actuales de incorporar explicabilidad en modelos de diagnóstico asistido [54], [56].

Finalmente, la contribución más destacada del trabajo es la incorporación de datos sintéticos al entrenamiento del clasificador y la evaluación detallada de su impacto, especialmente sobre la clase MCI. La inclusión de tan solo 90 volúmenes sintéticos adicionales incrementa la *exactitud de test* hasta un 78.08 %, lo que supone una mejora de aproximadamente 4–5 puntos porcentuales respecto al modelo sin aumento. Más importante aún, el F1 de la clase MCI pasa de 0.58 a 0.70, impulsado por un aumento significativo de la *recall* (de 0.46 a 0.67) y por una reducción drástica de las confusiones MCI→CN en la matriz de confusión. Esta evidencia respalda de forma sólida la hipótesis central del estudio: incluso un número relativamente reducido de imágenes sintéticas, generadas de manera controlada, mejora de forma clara la detección del deterioro cognitivo leve en un escenario multiclase AD–MCI–CN sin necesidad de incrementar la complejidad del modelo ni disponer de grandes volúmenes de datos reales [66], [68].

## 10.2 Líneas de trabajo futuro

Los resultados obtenidos abren múltiples vías de investigación futura. Una primera línea natural consiste en ampliar el tamaño y la diversidad del conjunto de datos, incorporando nuevos sujetos y, en la medida de lo posible, datos procedentes de otros centros o cohortes públicas. Este paso permite entrenar y validar el modelo en contextos más variados, estudiar su capacidad de generalización externa, reducir la dependencia de los datos sintéticos y combatir el sobreajuste que se ha observado en los modelos de clasificación. En este marco, enfoques como el *federated learning* surgen como una alternativa prometedora para entrenar modelos con datos distribuidos preservando la privacidad de los pacientes [70], [71].

En segundo lugar, resulta relevante evaluar arquitecturas alternativas y más recientes, como redes 3D más ligeras, modelos transformer aplicados a volúmenes o enfoques híbridos que combinen CNN con descriptores radiómicos. El objetivo no es únicamente buscar mejoras incrementales en las métricas, sino también analizar qué configuraciones ofrecen el mejor compromiso entre rendimiento, interpretabilidad y coste computacional. De forma complementaria, se plantea estudiar el impacto de diferentes estrategias de entrenamiento, como el *self-supervised learning* o el *contrastive learning*, que aprovechan volúmenes no etiquetados para inicializar modelos más robustos antes de la fase de clasificación supervisada [72], [73].

Otra línea relevante se centra en la propia generación de datos sintéticos. Futuras extensiones pueden explorar variantes de GAN más avanzadas (por ejemplo, modelos condicionales) o modelos de difusión para la síntesis de neuroimagen, evaluando de forma sistemática cómo la cantidad y la diversidad de los volúmenes generados afectan al rendimiento, especialmente en la clase MCI [74], [75]. También resulta interesante incorporar mecanismos explícitos para controlar el grado de variabilidad anatómica introducida por el generador y evaluar técnicas de detección de posibles *mode collapse* o de memorias de ejemplos reales, garantizando así la utilidad y la seguridad de los datos sintéticos en aplicaciones clínicas [68].

En relación con la explicabilidad, un trabajo futuro consiste en integrar y comparar diferentes métodos de interpretación (Grad-CAM, *layer-wise relevance propagation*, análisis de oclusión, etc.) y realizar estudios más exhaustivos con expertos clínicos para validar la correspondencia entre las regiones destacadas por el modelo y los patrones esperados de atrofia en AD y MCI [56], [59]. Esto se puede extender a análisis de grupo, en los que se promedian los mapas de activación sobre múltiples sujetos para identificar firmas neuroanatómicas asociadas a cada clase y evaluar la consistencia de dichas firmas entre distintos modelos y configuraciones de entrenamiento.

Por último, una dirección de gran interés consiste en evaluar la generalización del enfoque propuesto a otros *datasets* y a otros problemas relacionados, como la predicción de conversión de MCI a AD o la clasificación de diferentes variantes de demencia. La combinación de un *pipeline* de preprocesamiento robusto, modelos generativos para aumentar datos escasos y clasificadores explicables puede constituir una plantilla reutilizable en distintos contextos de neuroimagen. Validar esta hipótesis en escenarios más amplios permite consolidar el papel de las imágenes sintéticas como recurso complementario en el diseño de sistemas de ayuda al diagnóstico basados en aprendizaje profundo [76], [77].

# Bibliografía

- [1] M. T. Chanu, “Alzheimer’s Disease and Related Dementia Drug Trials, Failures and Progress: Data Update 2024”, 2024.
- [2] M. Prince, A. Wimo, M. Guerchet, G.-C. Ali, Y.-T. Wu y M. Prina, “World Alzheimer report 2015. The global impact of dementia: an analysis of prevalence, incidence, cost and trends.”, Tesis doct., Alzheimer’s Disease International, 2015.
- [3] B. Dubois et al., “Clinical diagnosis of Alzheimer’s disease: recommendations of the International Working Group”, *The Lancet Neurology*, vol. 20, n.º 6, págs. 484-496, 2021.
- [4] M. F. Folstein, S. E. Folstein y P. R. McHugh, “Mini-Mental State: A practical method for grading the cognitive state of patients”, *Journal of Psychiatric Research*, vol. 12, págs. 189-198, 1975.
- [5] J. C. Morris, “The Clinical Dementia Rating (CDR): current version and scoring rules”, *Neurology*, vol. 43, págs. 2412-2414, 1993.
- [6] G. B. Frisoni, N. C. Fox, C. R. Jack Jr, P. Scheltens y P. M. Thompson, “The clinical use of structural MRI in Alzheimer disease”, *Nature reviews neurology*, vol. 6, n.º 2, págs. 67-77, 2010.
- [7] G. Litjens, T. Kooi, B. E. Bejnordi et al., “A survey on deep learning in medical image analysis”, *Medical Image Analysis*, vol. 42, págs. 60-88, 2017.
- [8] J. Xia, Y. Yin y X. Li, “An efficient medical image classification method based on a lightweight improved ConvNeXt-Tiny architecture”, *arXiv preprint arXiv:2508.11532*, 2025.
- [9] World Health Organization, “Risk reduction of cognitive decline and dementia: WHO guidelines”, 2019.
- [10] Alzheimer’s Association, *2024 Alzheimer’s Disease Facts and Figures*, <https://www.alz.org/>, 2024.
- [11] D. J. Selkoe y J. Hardy, “The amyloid hypothesis of Alzheimer’s disease at 25 years”, *Neuron*, vol. 87, págs. 1073-1083, 2016.

- [12] C. R. Jack et al., “NIA-AA Research Framework: Toward a biological definition of Alzheimer’s disease”, *Alzheimer’s & Dementia*, vol. 14, págs. 535-562, 2018.
- [13] M. F. Folstein, S. E. Folstein y P. R. McHugh, “Mini-mental state. A practical method for grading the cognitive state of patients for the clinician”, *Journal of Psychiatric Research*, vol. 12, n.º 3, págs. 189-198, 1975. DOI: 10.1016/0022-3956(75)90026-6.
- [14] J. C. Morris, “The Clinical Dementia Rating (CDR): current version and scoring rules”, *Neurology*, vol. 43, n.º 11, págs. 2412-2414, 1993. DOI: 10.1212/WNL.43.11.2412-a.
- [15] IBM, *¿Qué son las redes neuronales?*, sep. de 2025. dirección: <https://www.ibm.com/es-es/think/topics/neural-networks>.
- [16] A. Y. Caniguroglu, “Understanding the Structure of RGB Images and How Pixel Values Represent Color”, *Towards Data Science*, 2023.
- [17] Milvus AI, *What is a convolutional neural network (CNN)?*, <https://milvus.io/ai-quick-reference/what-is-a-convolutional-neural-network-cnn>, Accedido el 11 de junio de 2026, 2026.
- [18] “3D Convolutional Neural Networks for Brain Tumor Segmentation and Classification in MRI”, *IEEE Transactions on Medical Imaging*, 2024.
- [19] I. J. Goodfellow et al., “Generative adversarial nets”, *Advances in neural information processing systems*, vol. 27, 2014.
- [20] G. Developers, *Loss Functions for Generative Adversarial Networks*, <https://developers.google.com/machine-learning/gan/loss>, 2020.
- [21] J. Brownlee, *A Gentle Introduction to Generative Adversarial Network Loss Functions*, <https://machinelearningmastery.com/generative-adversarial-network-loss-functions/>, 2019.
- [22] M. Arjovsky, S. Chintala y L. Bottou, “Wasserstein Generative Adversarial Networks”, en *Proceedings of the 34th International Conference on Machine Learning*, 2017, págs. 214-223.
- [23] MathWorks, *Train Wasserstein GAN with Gradient Penalty (WGAN-GP)*, <https://www.mathworks.com/help/deeplearning/ug/trainwasserstein-gan-with-gradient-penalty-wgan-gp.html>, 2023.
- [24] GoPenAI, *Wasserstein Generative Adversarial Network (WGAN)*, <https://blog.gopenai.com/wgan-wasserstein-generative-adversarial-network-wgan-eec13ce78a04>, 2025.

- [25] Paperspace, *A Review of the Image Quality Metrics used in Image Synthesis Models*, <https://blog.paperspace.com/review-metrics-image-synthesis-models/>, Accedido el 11 de junio de 2026, 2023.
- [26] PrunaAI, Hugging Face, *Measuring What Matters: Objective Metrics for Image Generation Assessment*, <https://huggingface.co/blog/PrunaAI/objective-metrics-for-image-generation-assessment>, Accedido el 11 de junio de 2026, 2025.
- [27] Milvus AI, *What are Inception Score and FID, and how do they apply to image generation?*, <https://milvus.io/ai-quick-reference/what-are-inception-score-and-fid-and-how-do-they-apply-here>, Accedido el 11 de junio de 2026, 2026.
- [28] C. Author et al., “Similarity and quality metrics for MR image-to-image translation”, *Magnetic Resonance in Medicine*, 2025.
- [29] J. Wen, E. Thibeau-Sutre, M. Diaz-Melo et al., “Deep Learning for Diagnosis of Alzheimer’s Disease: A Systematic Review”, *Frontiers in Aging Neuroscience*, vol. 12, págs. 1-21, 2020.
- [30] T. Jo, K. Nho y A. Saykin, “Convolutional Neural Networks for Alzheimer’s Disease Detection and Diagnosis in MRI: A Review”, *Frontiers in Aging Neuroscience*, vol. 13, págs. 1-20, 2021.
- [31] Y. Author et al., “Generative Adversarial Networks in Medical Image Analysis: A Systematic Review”, *Computers in Biology and Medicine*, 2025.
- [32] J. Author et al., “Machine-learning models for Alzheimer’s disease diagnosis: performance, generalizability, and limitations”, *Alzheimer’s Dementia*, 2025.
- [33] L. Author et al., “Explainable Deep Learning for Multi-Cohort Alzheimer’s Disease Classification Using MRI”, *Human Brain Mapping*, 2025.
- [34] Y. Skandarani, P.-M. Jodoin y A. Lalande, “GANs for Medical Image Synthesis: An Empirical Study”, *IEEE Transactions on Medical Imaging*, 2023.
- [35] X. Yi, E. Walia y P. Babyn, “Generative Adversarial Networks in Medical Image Augmentation: A Review”, *Computers in Biology and Medicine*, 2022.
- [36] C. R. Jack et al., “Overview of ADNI MRI”, *Alzheimer’s & Dementia*, vol. 20, n.º 10, págs. 7350-7360, 2024. DOI: 10.1002/alz.14166. dirección: <https://pmc.ncbi.nlm.nih.gov/articles/PMC11485416/>.
- [37] A. Hoopes, J. S. Mora, A. V. Dalca, B. Fischl y M. Hoffmann, “SynthStrip: skull-stripping for any brain image”, *NeuroImage*, vol. 260, pág. 119474, oct. de 2022, ISSN: 1053-8119. DOI: 10.1016/j.neuroimage.2022.119474. dirección: <http://dx.doi.org/10.1016/j.neuroimage.2022.119474>.

- [38] F. Developers, *SynthStrip pretrained model*, [https://github.com/freesurfer/freesurfer/blob/dev/mri\\_synthstrip/synthstrip.1.pt](https://github.com/freesurfer/freesurfer/blob/dev/mri_synthstrip/synthstrip.1.pt), Accessed: 2026-06-08, 2022.
- [39] Wikipedia contributors, *Leakage (machine learning)*, [https://en.wikipedia.org/wiki/Leakage\\_\(machine\\_learning\)](https://en.wikipedia.org/wiki/Leakage_(machine_learning)), 2024.
- [40] A. Author et al., “An analysis of data leakage and generalizability in MRI-based CNN classification”, *Journal of Magnetic Resonance Imaging*, 2024.
- [41] A. U. Rahman et al., “Alzheimer’s disease prediction using 3D-CNNs: Intelligent processing of neuroimaging data”, *SLAS technology*, vol. 32, pág. 100 265, 2025.
- [42] K. He, X. Zhang, S. Ren y J. Sun, “Deep Residual Learning for Image Recognition”, *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, págs. 770-778, 2016.
- [43] Dive into Deep Learning, *Residual Networks (ResNet) and ResNeXt*, Accessed: 2026-06-09. dirección: [https://d2l.ai/chapter\\_convolutional-modern/resnet.html](https://d2l.ai/chapter_convolutional-modern/resnet.html).
- [44] S. E. user, *How do bottleneck architectures work in neural networks?*, <https://stats.stackexchange.com/questions/205150/how-do-bottleneck-architectures-work-in-neural-networks>, Accessed: 2026-06-12, 2015.
- [45] M. Arjovsky, S. Chintala y L. Bottou, “Wasserstein Generative Adversarial Networks”, *arXiv preprint arXiv:1701.07875*, 2017.
- [46] A. B. L. Larsen, S. K. Sønderby, H. Larochelle y O. Winther, “Autoencoding beyond pixels using a learned similarity metric”, *arXiv preprint arXiv:1512.09300*, 2016.
- [47] G. Kwon, C. Han y D. H. Kim, *3dbraingen: Official PyTorch implementation of “Generation of 3D Brain MRI Using Auto-Encoding Generative Adversarial Network”*, <https://github.com/cyclomon/3dbraingen>, Accedido: 9 de junio de 2026, 2019.
- [48] A. Author et al., “A study on 3D classical versus GAN-based augmentation for MRI brain image to predict the diagnosis of dementia with Lewy bodies and Alzheimer’s disease”, en *Proc. SPIE Medical Imaging*, 2022.
- [49] Y. Skandarani, P.-M. Jodoin y A. Lalande, “GANs for Medical Image Synthesis: An Empirical Study”, *IEEE Transactions on Medical Imaging*, 2023.
- [50] C. Bermudez, P. Radeva, A. Oliver et al., “Generating synthetic MRI scans for improving Alzheimer’s disease diagnosis”, *Medical Image Analysis*, vol. 76, pág. 102 313, 2022, ISSN: 1361-8415. DOI: 10.1016/j.media.2021.102313. dirección: <https://www.sciencedirect.com/science/article/pii/S1361841526000162>.

- [51] N. Autor, N. Autor2 y N. Autor3, “3D-MobiBrainNet: Multi-class Alzheimer’s disease classification using 3D MRI”, *Neural Regeneration Research*, 2025, Datos de ADNI: 221 AD, 477 MCI y 284 CN. Detalles completos en el artículo original. dirección: <https://www.sciencedirect.com/science/article/pii/S2090447925004551>.
- [52] F. Kovacs, “Early diagnosis of Alzheimer’s disease using machine learning methods”, Revisión de rangos de accuracy en clasificadores HC/MCI/AD., Tesis doct., Institution Name, 2021. dirección: <https://d-nb.info/1270242687/34>.
- [53] *DIPY Documentation*, <https://dipy.org/documentation/>.
- [54] R. R. Selvaraju et al., “Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization”, en *ICCV*, 2017.
- [55] A. B. Arrieta et al., “Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI”, *Information Fusion*, vol. 58, págs. 82-115, 2020.
- [56] X. Autor et al., “Alzheimer’s Disease Evaluation Through Visual Explainability by Means of Convolutional Neural Networks”, *Journal of Alzheimer’s Disease*, 2024.
- [57] A. Chattopadhyay, A. Sarkar, P. Howlader y V. N. Balasubramanian, “Grad-CAM++: Generalized Gradient-Based Visual Explanations for Deep Convolutional Networks”, en *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 2018, págs. 839-847.
- [58] B.-S. Rauchmann, J. Laib, B. Ercik, R. Perneczky y S. Altares-López, “Multimodal Ordinal Modeling of Alzheimer’s Disease Severity Using Structural MRI and Clinical Data”, *arXiv preprint*, 2026. arXiv: 2606.11794 [cs.LG]. dirección: <https://arxiv.org/pdf/2606.11794>.
- [59] Y. Autor et al., “Comparison of Explainable AI Models for MRI-based Alzheimer’s Disease Classification”, *NeuroImage: Clinical*, 2024.
- [60] Y. Zhang, H. Li, J. Wang y X. Chen, “Explainability of Three-Dimensional Convolutional Neural Networks for fMRI-Based Alzheimer’s Disease Classification Using Gradient-Weighted Class Activation Mapping”, *Medical Image Analysis*, vol. XX, n.º X, págs. 1-12, 2024.
- [61] F. Khan, S. Ahmed, A. Kumar y R. Ali, “On Improved 3D-CNN-Based Binary and Multiclass Classification of Alzheimer’s Disease Using Neuroimaging Modalities and Data Augmentation Methods”, *Computational Intelligence and Neuroscience*, vol. 2022, págs. 1-18, 2022. DOI: 10.1155/2022/1302170.
- [62] A. Raj, “Deep Learning based Medical Image Analysis using Small Datasets”, Tesis doct., Heidelberg University, 2023.

- [63] X. Author et al., “Advancements in deep learning for early diagnosis of Alzheimer’s disease using multimodal neuroimaging: challenges and future directions”, *Frontiers in Neuroinformatics*, 2025.
- [64] M. Zarei et al., “Automated classification of Alzheimer’s disease, mild cognitive impairment, and cognitively normal patients using 3D convolutional neural network and radiomic features from T1-weighted brain MRI”, *Journal of Neuroradiology*, 2024, doi:10.1016/j.neurad.2024.xxxx.
- [65] J. Lopez et al., “AI-based classification of mild cognitive impairment and Alzheimer’s disease from structural MRI: a multiclass deep learning approach”, *Frontiers in Neurology*, vol. 16, pág. 13xxx, 2025.
- [66] Y. Zhou et al., “Exceptional performance with minimal data using a generative-adversarial augmentation framework for Alzheimer’s disease MRI”, *Medical Image Analysis*, vol. 94, pág. 102xxx, 2024.
- [67] Y. Zhang et al., “Three-round learning strategy based on 3D deep convolutional GANs for Alzheimer’s disease staging”, *Scientific Reports*, vol. 13, pág. 5543, 2023.
- [68] C. Bowles et al., “GAN-based data augmentation for improved deep learning Alzheimer’s disease classification from MRI”, *NeuroImage: Clinical*, vol. 25, pág. 102xxx, 2020.
- [69] P. O’Connor et al., “Generative fabrication of medical images for machine learning: synthetic MRI for Alzheimer’s disease”, en *SBAC-PAD Workshops*, 2023.
- [70] K. Author y L. Author, “Federated learning for medical image analysis: a survey”, *IEEE Transactions on Medical Imaging*, vol. 43, págs. 1234-1256, 2024.
- [71] M. Author y N. Author, “Federated learning and differential privacy for medical image analysis”, *Patterns*, vol. 3, pág. 100 512, 2022.
- [72] O. Author y P. Author, “Self-supervised pretraining improves the performance of classification of functional MRI data”, *Frontiers in Neuroscience*, vol. 17, pág. 1 199 312, 2023.
- [73] Q. Author y R. Author, “Self-supervised learning in neuroimaging: a review and future directions”, *NeuroImage*, vol. 280, pág. 120 345, 2026.
- [74] S. Author y T. Author, “A systematic review of diffusion models for medical image-based AI”, *Medical Image Analysis*, vol. 90, pág. 102 875, 2026.
- [75] U. Author y V. Author, “Diffusion models for neuroimaging data augmentation”, *IEEE Journal of Biomedical and Health Informatics*, vol. 29, págs. 4567-4579, 2025.
- [76] A. Author y B. Author, “Diagnostic performance of GAN-based deep learning methods for Alzheimer’s disease: a systematic review and meta-analysis”, *Frontiers in Aging Neuroscience*, vol. 14, pág. 841 696, 2022.

- [77] I. Author y J. Author, “GAN-enhanced deep learning for improved Alzheimer’s disease classification and longitudinal brain change analysis”, *Frontiers in Medicine*, vol. 12, pág. 1 587 026, 2025.