

UNIVERSIDAD COMPLUTENSE DE MADRID
FACULTAD DE CIENCIAS MATEMÁTICAS

UNIVERSIDAD POLITÉCNICA DE MADRID
ESCUELA TÉCNICA SUPERIOR DE INGENIEROS DE
TELECOMUNICACIÓN

**MÁSTER EN TRATAMIENTO ESTADÍSTICO
COMPUTACIONAL DE LA INFORMACIÓN**



TRABAJO DE FIN DE MÁSTER

Redes Generativas Antagónicas

Antón Makarov Samusev

Director

Francisco Javier Yáñez Gestoso

Madrid, 2019

Resumen

En este trabajo estudiamos en profundidad las Redes Generativas Antagónicas, un modelo generativo no supervisado que explora la idea de poner a competir entre sí a dos redes neuronales. Expandimos las demostraciones teóricas del artículo en las que fueron introducidas por primera vez [7] y proponemos la utilización de una variante de este tipo de redes, llamada DCGAN [16], para la generación de imágenes de cuadros de arte, demostrando experimentalmente que es posible obtener imágenes visualmente atractivas mediante este método.

Abstract

In this dissertation we study Generative Adversarial Networks, an unsupervised generative model that explores the idea of setting two neural networks to compete with each other. We expand the proofs from the paper where they were first introduced [7] and propose the use of a variant of these kind of networks called DCGAN [16] to generate images of art, proving experimentally that it is possible to obtain visually appealing images with this method.

Índice general

1. Introducción	7
2. Preliminares	10
2.1. Aprendizaje automático	10
2.1.1. Aprendizaje supervisado	10
2.1.2. Aprendizaje no supervisado	11
2.1.3. Aprendizaje por refuerzo	11
2.2. Conceptos básicos	11
2.3. Redes neuronales	12
2.4. Aprendizaje profundo	12
2.4.1. Redes neuronales convolucionales	12
3. Redes generativas antagónicas	18
3.1. Idea general	18
3.2. Bases teóricas	19
3.2.1. Óptimo global	19
3.2.2. Convergencia	22
4. Generación de arte	24
4.1. DCGAN	24
4.1.1. Preprocesado	26
4.1.2. Arquitectura	27
4.1.3. Recursos	28
4.2. Resultados	28
5. Conclusión	31
Bibliografía	32

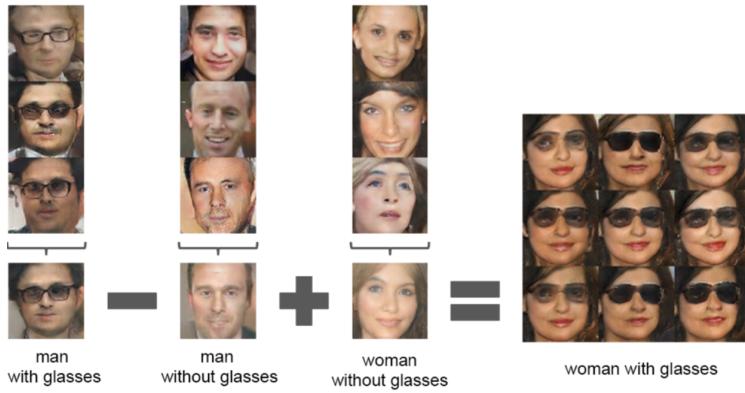
Capítulo 1

Introducción

El aprendizaje automático es uno de los temas que más interés está suscitando en los últimos años, tanto por parte de la comunidad investigadora como por la industria. Podemos justificar esta tendencia mediante la disponibilidad de recursos computacionales con costes ajustados, permitiendo la utilización en la práctica de algunos algoritmos antiguos de manera eficaz y promoviendo la creación de otros nuevos. También podemos asociar este éxito a la creciente disponibilidad de datos, que cada vez son recogidos por más empresas y organismos públicos de manera electrónica, haciendo más sencillo su manejo.

En este trabajo vamos a centrar nuestro estudio en las Redes Generativas Antagónicas (**GANs** o Generative Adversarial Networks en inglés), un tipo de modelo generativo que está recibiendo mucha atención en la actualidad debido a la introducción de ideas innovadoras en el campo del aprendizaje automático y a la cantidad de aplicaciones que posee. Los modelos generativos tienen como objetivo describir la distribución subyacente de los conjuntos de datos con los que trabajan, por ejemplo, podrían aproximar la distribución de los valores matemáticos que caracterizan las imágenes de determinados cuadros de arte para luego tomar muestras y, de este modo, crear cuadros totalmente nuevos de apariencia realista. Las GANs fueron propuestas por Ian Goodfellow en 2014 [7], donde describe un nuevo método de entrenamiento para modelos generativos basado en el concepto de equilibrio de teoría de juegos. A grandes rasgos, y en el caso del análisis de imágenes de cuadros de arte, propone un juego en el que dos redes neuronales compiten entre sí, una generando imágenes a partir de ruido aleatorio y la otra decidiendo si son reales o falsas, dando lugar a un entrenamiento en bucle, en el que ambas redes aprenden la una de la otra.

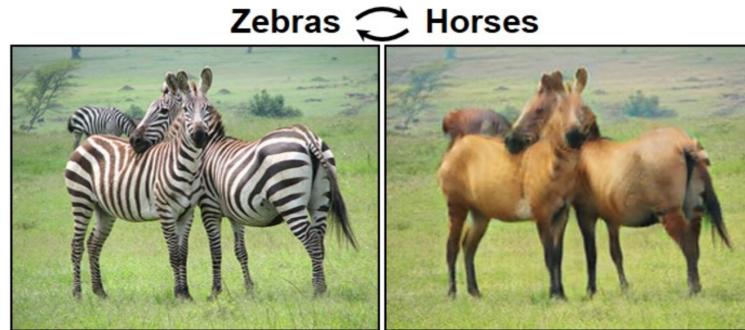
La acogida de esta idea ha sido espectacular, llevando a la publicación de numerosos artículos con extensiones, mejoras y modificaciones, obteniendo resultados impresionantes. Entre dichos artículos consideramos que merecen una mención especial *DCGAN* [16], que introduce una gran cantidad de mejoras al artículo de Goodfellow y describe el espacio latente de las GANs (Figura 1.1a), *Progressive Growing of GANs* [10], que permite generar imágenes de alta resolución mediante el entrenamiento progresivo de las redes (Figura 1.1b), *Cycle-GAN* [20], que permite la traducción de imágenes según su estilo, conservando el contexto (Figura 1.1c) y *Wasserstein GAN* [1] que describe una variante de GAN con muy buenas propiedades matemáticas.



(a) DCGAN



(b) Progressive growing of GANs



(c) Cycle-GAN

Figura 1.1: Imágenes descriptivas de algunas variantes de GANs.

En el Capítulo 2 introduciremos los conceptos esenciales del aprendizaje automático, daremos algunas definiciones que serán de utilidad en capítulos posteriores y describiremos brevemente las redes neuronales, centrando la atención en las redes neuronales convolucionales. En el Capítulo 3 explicaremos en mayor detalle la idea conceptual de las redes generativas antagónicas y desarrollaremos la teoría matemática en la que se fundamentan. El Capítulo 4 está dedicado a la descripción de nuestra implementación en Python de una DCGAN para generar imágenes artísticas. Este capítulo va ligado al repositorio de código con el que acompañamos el proyecto¹, que constituye una parte fundamental del trabajo realizado. Finalmente, en el Capítulo 5,

¹Todo el código de este proyecto se puede encontrar en el repositorio de GitHub <https://github.com/ant-mak/tfm>, donde además incluimos las instrucciones para la reproducción de los resultados del proyecto.

presentamos las conclusiones y señalamos algunos temas y líneas de investigación que sería interesante perseguir en el futuro.

Capítulo 2

Preliminares

En este capítulo realizaremos un breve resumen de los conceptos fundamentales del aprendizaje automático, introduciremos las definiciones necesarias para el desarrollo del resto del trabajo y daremos unas nociones básicas sobre aprendizaje profundo, prestando especial atención a las redes convolucionales.

2.1. Aprendizaje automático

El aprendizaje automático (**ML** o Machine Learning en inglés) tiene como objetivo principal el desarrollo de técnicas que permitan aprender a los ordenadores de manera autónoma, sin ser programadas explícitamente para ello, a partir de datos. Dentro del aprendizaje automático existen problemas de muy diversa índole; una manera de clasificarlos es la basada en el tipo de información disponible. En las siguientes secciones describiremos brevemente la clasificación más habitual.

2.1.1. Aprendizaje supervisado

Los problemas de aprendizaje supervisado se caracterizan por la disponibilidad tanto de una serie de muestras X como de la información sobre cuál debe ser el resultado y para cada una de ellas. Buscan encontrar una función que, a partir de una muestra, devuelva el resultado y correspondiente. Algunos de los problemas típicos ante los que nos podemos encontrar dentro del aprendizaje supervisado son:

- Predecir el precio de un piso a partir del número de habitaciones, de los metros cuadrados que tiene, del barrio, etc.
- Predecir si una transacción es o no fraudulenta a partir de la cantidad transferida, del lugar en el que ha sido realizada, de las transacciones realizadas durante los días previos, etc.
- Determinar si una imagen es un perro o un gato.

Algunos de los algoritmos más conocidos para la resolución de problemas de aprendizaje supervisado son la regresión lineal y logística, las máquinas de vector soporte, los árboles de decisión o las redes neuronales.

2.1.2. Aprendizaje no supervisado

En el caso de los problemas de aprendizaje no supervisado tan solo se dispone de las muestras X , sin etiqueta alguna, y se busca recuperar la estructura de los datos. Algunos de los problemas típicos del aprendizaje no supervisado son:

- Segmentación para campañas de publicidad, en las que se dispone de datos socioeconómicos de clientes y se les quiere dividir en varias clases según sus intereses.
- Generación de cuadros nuevos a partir de imágenes de cuadros reales.

Algunos algoritmos utilizados en aprendizaje no supervisado son el K-means o las redes neuronales.

2.1.3. Aprendizaje por refuerzo

El aprendizaje por refuerzo busca determinar la mejor acción a realizar ante una situación concreta para maximizar la recompensa obtenida por dicha acción. Algunos problemas dentro de este tipo de aprendizaje son:

- Enseñan a una máquina a jugar al ajedrez, Go, StarCraft, etc.
- Enseñar a una máquina a conducir.

Los algoritmos más conocidos dentro de este área son el Q-learning, SARSA o DQN.

2.2. Conceptos básicos

Dedicaremos esta sección a la definición de algunos conceptos que serán de utilidad en capítulos posteriores.

Definición 2.1. *La divergencia de Kullback-Leibler es una medida de la diferencia entre dos distribuciones. Se define como:*

$$D_{\text{KL}}(P\|Q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx. \quad (2.1)$$

Definición 2.2. *La divergencia de Jensen-Shannon se define a partir de la de Kullback-Leibler como:*

$$D_{\text{JS}}(P\|Q) = \frac{1}{2}D_{\text{KL}}(P\|M) + \frac{1}{2}D_{\text{KL}}(Q\|M), \quad (2.2)$$

donde $M = \frac{P+Q}{2}$. Tiene la ventaja de ser simétrica y finita.

Definición 2.3. *Un equilibrio de Nash es un concepto de solución para juegos no cooperativos de dos o más jugadores en el que cada jugador conoce las estrategias de equilibrio de los demás. Si cada jugador elige una estrategia y ninguno de ellos tiene incentivos para cambiar la suya suponiendo que los demás no lo hacen, se dice que el conjunto de estrategias está en equilibrio.*

Finalmente, introducimos la **ley del inconsciente estadístico**, utilizada para calcular la esperanza de una función $f(X)$ de una variable aleatoria X cuando se conoce la distribución de X pero no la de $f(X)$:

$$\mathbb{E}_{x \sim p}[f(x)] = \int p(x)f(x) dx. \quad (2.3)$$

2.3. Redes neuronales

Las redes neuronales son una clase de algoritmos de aprendizaje automático inspiradas en el estudio biológico de las neuronas en el cerebro. Cabe destacar que pese a que, toman la idea de las neuronas biológicas, no pretenden modelizarlas ni replicar su funcionamiento. Están formadas por una gran cantidad de unidades de procesamiento (o neuronas) interconectadas entre sí. Las fuerzas (o pesos) de cada una de dichas conexiones se modifica en el proceso de entrenamiento con el objetivo de minimizar una cierta función de coste. La minimización se lleva a cabo mediante el algoritmo de descenso en la dirección del gradiente o alguna de sus variantes creadas específicamente para las redes neuronales (SGD, RMSprop, Adam, etc.). Para calcular el gradiente se utiliza el algoritmo de retropropagación, basado en la regla de la cadena. Existe una gran variedad de redes neuronales según el tipo de problema de aprendizaje que se quiera resolver con ellas, la caracterización de las neuronas y la estructura de conexiones de la red. Algunas de las más conocidas son el perceptrón multicapa, las redes convolucionales o las redes recurrentes.

2.4. Aprendizaje profundo

El aprendizaje profundo es un subgrupo de algoritmos de aprendizaje automático que buscan abstraer características de alto nivel presentes en los datos, utilizando generalmente arquitecturas basadas en redes neuronales con una gran cantidad de capas. Para profundizar en las matemáticas detrás del aprendizaje profundo una excelente referencia es [6].

En la actualidad son muy populares debido a su enorme potencial para resolver problemas complejos en ámbitos como la visión por computador o el procesamiento del lenguaje natural, entre otros. Dicho potencial no ha podido ser aprovechado hasta hace muy poco debido a la escasez, por un lado, de conjuntos de datos adecuados y, por otro, de recursos computacionales, ya que estos métodos requieren unas capacidades de cálculo y dimensiones de conjuntos de datos muy superiores a otros métodos más tradicionales.

2.4.1. Redes neuronales convolucionales

Dado que nuestro objetivo es describir las redes generativas antagónicas y, más concretamente, utilizarlas para la generación de imágenes, es imprescindible introducir algunos de los conceptos más importantes de la familia de arquitecturas de aprendizaje profundo que están a la vanguardia en el campo de la visión por computador, las redes neuronales convolucionales (**CNNs** o Convolutional Neural Networks en inglés).

Debemos sus primeros desarrollos a Kunihiko Fukushima, que en 1980 introdujo el Neocognitrón [4] que, posteriormente, sería tomado por Yann LeCun [13], que describió la mayoría de conceptos tal como los conocemos hoy en día.

La particularidad de las CNNs es que introducen una visión local del espacio. Pensemos, por ejemplo, en una imagen de un perro. Parece razonable que intentemos identificar que efectivamente nos encontramos ante un perro y no un gato fijándonos en pequeñas partes de la imagen como ojos, hocico, patas, orejas, etc., en lugar de en todos los píxeles a la vez sin ningún tipo de relación espacial entre sí. Es destacable que este tipo de filtros ya se utilizaba antes de las redes convolucionales, detectando por ejemplo líneas horizontales o verticales y seleccionando regiones de interés manualmente. Sin embargo, las redes convolucionales van más allá, automatizando el proceso de extracción de características y aprendiendo durante el entrenamiento los filtros más idóneos para el problema.

Veamos ahora en detalle el funcionamiento de las redes convolucionales a la vez que desarrollamos los conceptos básicos que nos serán de utilidad a lo largo del resto del trabajo.

Definición 2.4. *Una imagen a color se puede definir como un tensor, con altura, anchura y profundidad. Las dos primeras dimensiones nos indican el tamaño de la imagen, por ejemplo 64×64 píxeles. La profundidad contiene los canales de color, generalmente rojo, verde y azul (**RGB** o Red Green Blue en inglés) con una cierta intensidad representada en un rango de 0 a 255. Mediante la combinación de dichos canales se forman las imágenes a color a las que estamos acostumbrados. En la Figura 2.1 tenemos una representación de esta idea.*

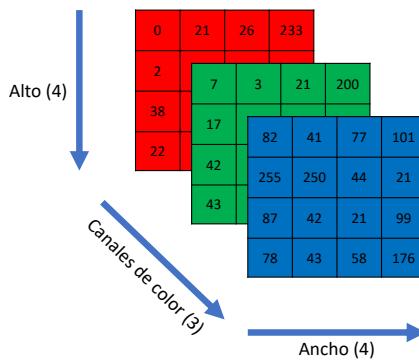


Figura 2.1: Representación esquemática de una imagen.

Una vez definida la imagen formalmente, podemos comenzar con la descripción de la capa de convolución. La función que desempeña es la de extraer características de una imagen de manera automática. Para ello, recorre

el tensor de entrada secuencialmente con una cuadrícula de tamaño pre-determinado, aplicando a cada una de las posiciones de la cuadrícula uno o varios filtros (**Kernels** en inglés) cuyos coeficientes son aprendidos por la red. La salida de una operación de convolución recibe el nombre de mapa de características.

En la Figura 2.2 tenemos una representación esquemática de todos los elementos involucrados. En la parte superior está representada una imagen de 5×5 píxeles en 3 canales: rojo, verde y azul. Los bordes de color gris son un relleno (**Padding** en inglés) de ceros cuya función es recoger mejor la información de los bordes de la imagen y modular las dimensiones de salida. La imagen está siendo recorrida por una cuadrícula 3×3 con un desplazamiento (**Stride** en inglés) de tamaño 2, es decir, en cada paso se mueve dos posiciones hacia la derecha y al llegar al final de la línea baja dos posiciones, colocándose de nuevo a la izquierda. En la parte inferior se sitúan los filtros, que generan de el mapa de características representado a la derecha. Estos filtros funcionan de la manera siguiente: primero se realiza un producto elemento a elemento entre la cuadrícula y el filtro en cada canal de color, se suman el resultado del producto y, finalmente, se suman los resultados de cada filtro, dando lugar a un elemento del mapa de características.

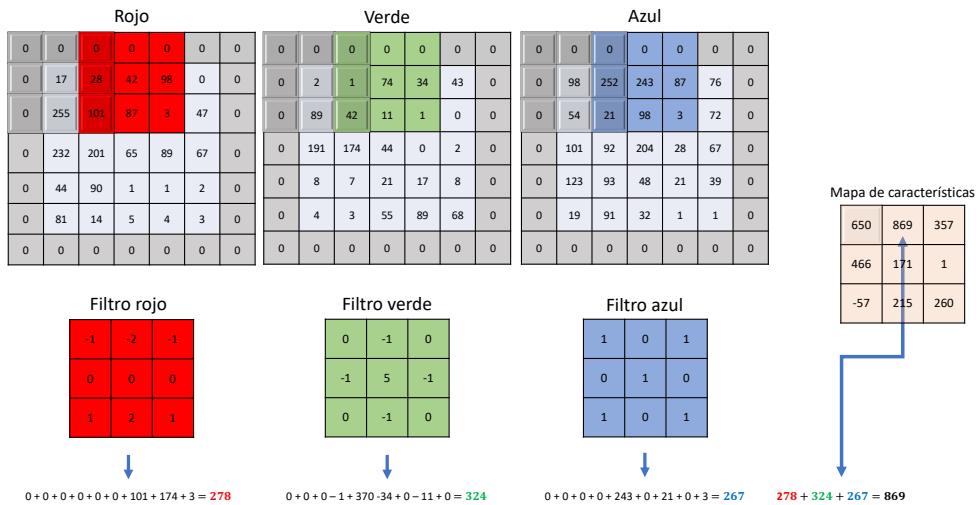


Figura 2.2: Diagrama ilustrativo de la operación de convolución.

En general se puede calcular la dimensión de salida de una capa convolucional mediante la siguiente fórmula:

$$O = \frac{W - K + 2P}{S} + 1,$$

donde W es el tamaño de entrada, K el tamaño del filtro, P el padding y S la longitud del paso. En nuestro ejemplo, por tanto, tendremos un tamaño de salida $O = \frac{5-3+2}{2} + 1 = 3$. En general se suele utilizar más de un filtro en cada capa convolucional, obteniendo así también dimensiones de profundidad.

Observación 2.1. *Como se puede ver en la fórmula, el parámetro stride influye en la dimensión de salida, cosa que utilizaremos extensivamente en las GANs.*

Después de la capa de convolución se suele realizar una reducción de muestreo (**Pooling** en inglés), con el objetivo de reducir la dimensión de los mapas de características y abstraer su contenido para conseguir que la red generalice mejor. Pensemos por ejemplo en un mapa de características que identifica líneas horizontales. No necesitamos saber exactamente dónde tenemos una línea horizontal, nos basta con saber que hay una en el centro de la imagen. La capa de pooling toma cada uno de los mapas de características obtenidos y lo divide en regiones. A los elementos de cada región les aplica una operación para comprimir la información que contienen. Las operaciones más utilizadas son el máximo, la media o el mínimo. En la Figura 2.3 observamos de manera esquemática un mapa de características 4×4 que se ha dividido en 4 regiones 2×2 , a las que se aplica la función máximo.

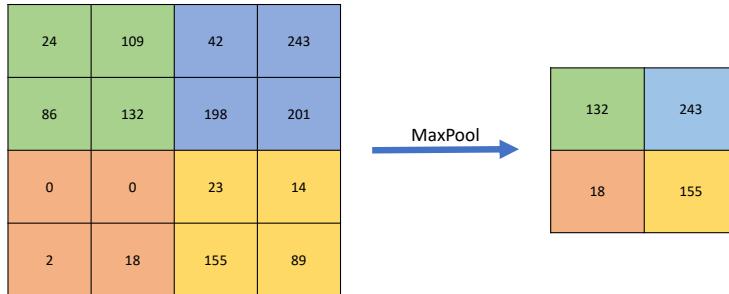


Figura 2.3: Diagrama ilustrativo de reducción de muestreo mediante máximo.

Antes de seguir adelante, es preciso introducir el concepto de **batch**, que, aunque aplica a todo tipo de redes neuronales, en nuestro caso tiene especial relevancia. El batch es un hiperparámetro que controla el número de muestras con el que es entrenada la red antes de proceder a una actualización de pesos. El tamaño del batch puede ir desde 1 hasta el número de muestras del conjunto de datos. Cuando se recorre todo el conjunto de datos se dice que ha pasado una **época**. Por ejemplo, en el caso en que el tamaño de batch sea 100 y el conjunto de datos tenga 1000 muestras, actualizaremos los pesos de la red cada 100 muestras y diremos que hemos completado una época cuando lleguemos a 1000.

Recientemente se ha incorporado una técnica para mejorar la estabilidad en el entrenamiento de redes convolucionales, llamada **Batch Normalization** [9]. La idea fundamental es, si normalizamos los datos antes de introducirlos en la capa de entrada para reducir las influencias de las magnitudes de las variables y ajustarlos a una distribución, ¿por qué no hacer lo mismo en todas las capas?

Dado que los algoritmos de optimización pueden deshacer la normalización que introduzcamos si con ello reducen el valor de la función objetivo, es

necesario introducir dos parámetros adicionales γ y β que serán aprendidos por la red. En el Algoritmo 1 mostramos un breve esquema del funcionamiento. Para más detalles se puede consultar [9].

Algoritmo 1: Batch Normalization.

- 1 Calcular la media del batch: $\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$;
 - 2 Calcular la varianza del batch: $\sigma_b^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$;
 - 3 Normalizar: $\hat{x}_i = \frac{x_i - \mu_b}{\sqrt{\sigma_B}}$;
 - 4 Escalar y trasladar: $y_i = \gamma x_i + \beta$;
-

Finalmente hablemos de las funciones de activación. Un problema común de las funciones tradicionales como la sigmoide o la tangente hiperbólica es que en redes profundas puede dar lugar al fenómeno conocido como desvanecimiento del gradiente, que ocurre cuando el gradiente se aproxima excesivamente a cero, impidiendo el entrenamiento de la red. Por ejemplo, en el caso de la tangente hiperbólica, tenemos que su derivada se encuentra entre cero y uno. Dado que las redes neuronales utilizan el algoritmo de retropropagación, que a su vez está basado en la regla de la cadena, estaríamos multiplicando cantidades menores que uno tantas veces como capas tengamos, reduciendo de este modo de manera considerable el gradiente y haciendo el uso de estas funciones en redes de muchas capas inapropiado.

Por este motivo es necesario definir otras funciones que sigan siendo no lineales pero tengan un mejor comportamiento ante estos escenarios.

Definición 2.5. La función de activación unidad lineal rectificada (**ReLU** o *REctified Linear Unit* en inglés) se define como:

$$f(x) = \max(0, x)$$

De manera similar, se define la función **LeakyReLU**:

$$f(x) = \begin{cases} x & \text{si } x > 0 \\ 0,01x & \text{en otro caso} \end{cases}$$

Observación 2.2. La principal ventaja de la función LeakyReLU es que evita que ciertas neuronas se apaguen indefinidamente, dejando pasar una pequeña parte del gradiente también para el caso en que $x < 0$.

Como mencionamos anteriormente, describir exhaustivamente las CNNs no es el objetivo de este trabajo, por lo que nos hemos dejado muchos elementos relevantes en el tintero. Aún así, ya estamos en condiciones de proporcionar la arquitectura básica de una red convolucional. Suele constar de dos partes; en la primera se busca extraer características de las imágenes mediante varias capas de convolución con pooling y activaciones ReLU, aumentando progresivamente la profundidad a la vez que se disminuye anchura y altura. Una vez construidos los mapas de características, se aplana y se procede a la parte de clasificación, en la que se puede utilizar por ejemplo un perceptrón multicapa. En la Figura 2.4 se puede observar una arquitectura de red convolucional como la que hemos descrito.

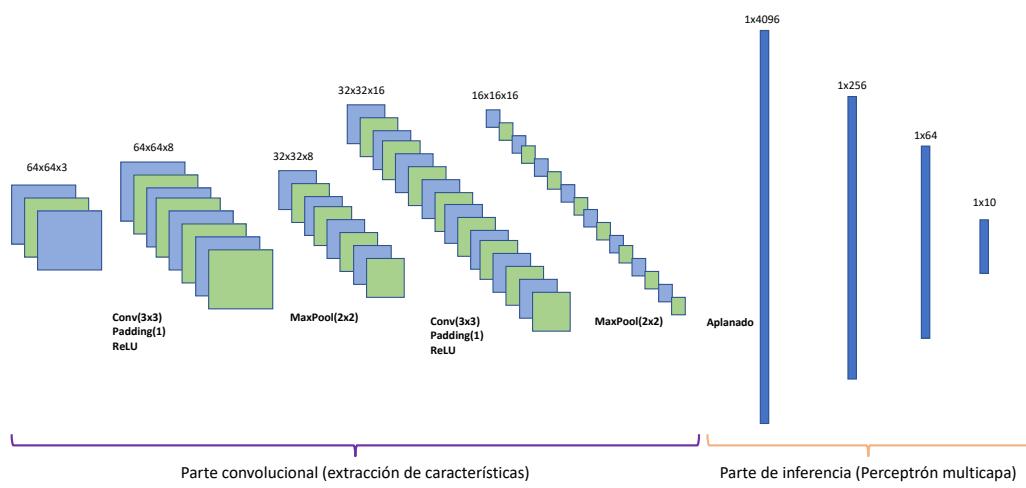


Figura 2.4: Arquitectura básica de una red convolucional

Capítulo 3

Redes generativas antagónicas

En este capítulo describiremos las ideas principales, tanto teóricas como prácticas de las redes generativas antagónicas. Daremos en primer lugar una idea conceptual de su funcionamiento, centrándonos después en la demostración de sus propiedades teóricas más relevantes. Finalmente concluimos con algunas consideraciones a tener en cuenta en la práctica.

3.1. Idea general

El objetivo de las GANs, y en general de los modelos generativos, es aprender la distribución que siguen los datos, pudiendo obtener así, en última instancia, muestras de dicha distribución. En general, las distribuciones que queremos modelar son muy complejas. Supongamos que nuestro objetivo es tomar muestras de la distribución de imágenes de perros, o dicho de otro modo, generar fotos de perros que sean realistas pero que no existan en la realidad ni sean una mezcla de imágenes de nuestro conjunto de entrenamiento. Tenemos la seguridad de que la distribución es extremadamente intrincada, existen perros de distintos colores, tamaños, razas, etc. Este problema es el que van a tratar de atacar las GANs.

La idea fundamental y más novedosa detrás de las GANs es poner dos redes neuronales a competir entre sí. Una red, llamada generadora (G), está dedicada a generar imágenes a partir de ruido aleatorio con distribución $p_z(z)$, mientras que otra, llamada discriminadora (D), trata de averiguar si la imagen es real o ficticia. Es frecuente ilustrar esta idea mediante la analogía de falsificadores de billetes que tratan de engañar a la policía. Los falsificadores empiezan dibujando billetes que no tienen nada que ver con los reales, intentando utilizarlos para realizar pagos, momento en el que son atrapados por la policía. Los falsificadores por tanto se dan cuenta de que están dibujando los billetes de manera incorrecta y modifican su técnica, mientras que la policía va aprendiendo a su vez a detectar mejor los billetes falsos. De este modo, a lo largo del entrenamiento se busca llegar a un equilibrio, en el que la policía no sea capaz de discernir los billetes falsos de los verdaderos, obteniendo así los ladrones una falsificación realista.

Traduzcamos ahora esta idea a términos matemáticos. Sea el generador G un PMC con parámetros θ_g que tiene por objetivo encontrar una distribución p_g lo más parecida posible a la distribución de los datos p_d . Para ello toma

un vector aleatorio z distribuido según una cierta distribución p_z y lo lleva al espacio definido por los datos. Por otro lado, sea el discriminador D otro PMC con parámetros θ_d , que, para una muestra, devuelve la probabilidad de que esta sea real (x) o falsa ($G(z)$). Entrenaremos D para que maximice la probabilidad de asignar la etiqueta correcta tanto a los ejemplos de entrenamiento como a los generados, mientras que, al mismo tiempo, entrenamos G para que minimice $\log(1 - D(G(z)))$. Es decir, tenemos el siguiente juego minimax con función de utilidad $V(G, D)$:

$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim p_d(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))] \quad (3.1)$$

En la Figura 3.1 mostramos una representación esquemática del funcionamiento de las GANs y en el Algoritmo 2 lo describimos de manera más concisa, dando además, una primera idea de cómo se lleva a cabo el entrenamiento.

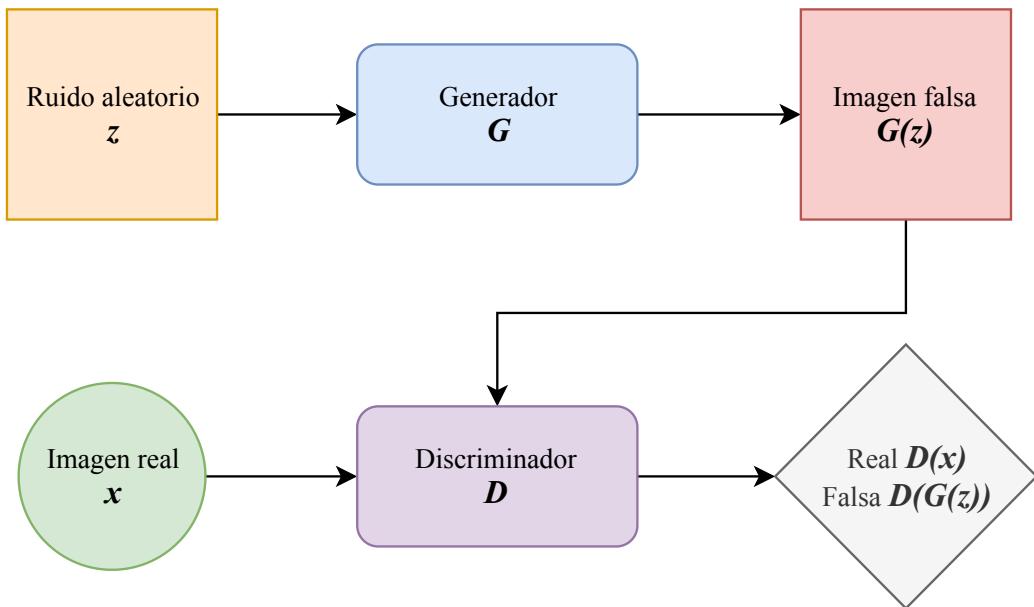


Figura 3.1: Diagrama conceptual de las redes generativas antagónicas.

3.2. Bases teóricas

En esta sección realizaremos un análisis de la teoría que hay detrás de las redes generativas antagónicas, utilizando como referencias principales [7, 5]. En primer lugar veremos cuáles son el generador y discriminador óptimos, mostrando después que el Algoritmo 2 logra que p_g converja a p_d .

3.2.1. Óptimo global

Veamos en primer lugar cual es el discriminador óptimo para un generador dado.

Teorema 3.1. *Para un generador G fijo, el discriminador D óptimo es:*

$$D_G^*(x) = \frac{p_d(x)}{p_d(x) + p_g(x)}. \quad (3.2)$$

Algoritmo 2: Entrenamiento minibatch de una red generativa antagonica. El valor de k es un hiperparámetro que se puede elegir libremente.

```

1 for Iteraciones de entrenamiento do
2   for  $k$  pasos do
3     Tomar muestra  $(z^{(1)}, z^{(2)}, \dots, z^{(m)})$  de tamaño  $m$  de  $p_g(z)$ ;
4     Tomar muestra  $(x^{(1)}, x^{(2)}, \dots, x^{(m)})$  de tamaño  $m$  de  $p_d(x)$ ;
5     Actualizar el discriminador ascendiendo su gradiente:
6       
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))].$$

7   end
8   Tomar muestra  $(z^{(1)}, z^{(2)}, \dots, z^{(m)})$  de tamaño  $m$  de  $p_g(z)$ ;
9   Actualizar el generador ascendiendo su gradiente:
10    
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

end

```

Demostración. El discriminador busca maximizar la Ecuación (3.1), que podemos reescribir, utilizando la ley del inconsciente estadístico como:

$$\begin{aligned} V(G, D) &= \int_x p_d(x) \log(D(x)) dx + \int_z p_g(z) \log(1 - D(g(z))) dz \\ &= \int_x (p_d(x) \log(D(x)) + p_g(x) \log(1 - D(x))) dx. \end{aligned}$$

Maximizando el integrando obtendremos el discriminador óptimo. Para ello, reescribimos el integrando como:

$$f(y) = a \log y + b \log(1 - y),$$

derivamos e igualamos a cero:

$$f'(y) = 0 \implies \frac{a}{y} - \frac{b}{1-y} = 0 \implies y = \frac{a}{a+b}.$$

Suponiendo que $a + b \neq 0$, hacemos la segunda derivada en $\frac{a}{a+b}$, obteniendo:

$$f''\left(\frac{a}{a+b}\right) = -\frac{a}{\left(\frac{a}{a+b}\right)^2} - \frac{b}{\left(1 - \frac{a}{a+b}\right)^2},$$

que es menor que cero para $(a,b) \in (0,1)$. Tomando por tanto $D = \frac{p_d}{p_d+p_g}$ tenemos el resultado, que además es único, ya que f tiene un solo máximo en el intervalo. \square

Es importante considerar que en la práctica, al desconocer $p_d(x)$, no podemos obtener el D óptimo, pero su existencia nos permitirá demostrar ahora que existe un óptimo para G .

Teorema 3.2. El mínimo global de $C(G) = \max_D V(G, D)$ se alcanza si y solo si $p_g = p_d$.

*Demuestra*ón. Para $p_g = p_d$, sustituyendo en la Ecuación (3.2), obtenemos que $D_G^* = \frac{1}{2}$. Podemos por tanto escribir:

$$\begin{aligned} V(G, D_G^*) &= \int_x \left(p_d(x) \log \frac{1}{2} + p_g(x) \log \left(1 - \frac{1}{2} \right) \right) dx \\ &= \int_x \log \frac{1}{2} (p_d(x) + p_g(x)) dx \\ &= -\log 2 \left(\int_x p_d(x) dx + \int_x p_g(x) dx \right) \\ &= -2 \log 2 = -\log 4. \end{aligned}$$

Así, nuestro valor candidato para ser el mínimo global es $-\log 4$. Veamos qué pasa si descartamos la hipótesis de que $p_g = p_d$. Para todo G , podemos sustituir D_G^* en $C(G)$:

$$C(G) = \int_x \left(p_d(x) \log \left(\frac{p_d(x)}{p_g(x) + p_d(x)} \right) + p_g(x) \log \left(\frac{p_g(x)}{p_g(x) + p_d(x)} \right) \right) dx. \quad (3.3)$$

Donde el segundo sumando de la integral proviene de:

$$\begin{aligned} 1 - D_G^*(x) &= 1 - \frac{p_d(x)}{p_g(x) + p_d(x)} \\ &= \frac{p_g(x) + p_d(x)}{p_g(x) + p_d(x)} - \frac{p_d(x)}{p_g(x) + p_d(x)} \\ &= \frac{p_g(x)}{p_g(x) + p_d(x)}. \end{aligned}$$

Sumando y restando $p_d(x) \log 2$ a ambos sumandos de (3.3) se tiene que:

$$\begin{aligned} C(G) &= \int_x \left((\log 2 - \log 2)p_d(x) + p_d(x) \log \left(\frac{p_d(x)}{p_g(x) + p_d(x)} \right) \right) dx \\ &\quad + \int_x \left((\log 2 - \log 2)p_g(x) + p_g(x) \log \left(\frac{p_g(x)}{p_g(x) + p_d(x)} \right) \right) dx. \end{aligned}$$

Podemos escribir esto como:

$$\begin{aligned} C(G) &= -\log 2 \int_x p_d(x) + p_g(x) dx \\ &\quad + \int_x p_d(x) \left(\log 2 + \log \left(\frac{p_d(x)}{p_g(x) + p_d(x)} \right) \right) dx \\ &\quad + \int_x p_g(x) \left(\log 2 + \log \left(\frac{p_g(x)}{p_g(x) + p_d(x)} \right) \right) dx. \end{aligned}$$

Simplificando términos obtenemos:

$$C(G) = -\log 4 + \int_x p_d(x) \log \left(\frac{p_d(x)}{\frac{p_g(x) + p_d(x)}{2}} \right) dx + \int_x p_g(x) \log \left(\frac{p_g(x)}{\frac{p_g(x) + p_d(x)}{2}} \right) dx. \quad (3.4)$$

Ahora podemos identificar la Ecuación (3.4) como la divergencia de Kullback-Leibler:

$$C(G) = -\log 4 + D_{\text{KL}} \left(p_d(x) \middle\| \frac{p_d(x) + p_g(x)}{2} \right) + D_{\text{KL}} \left(p_g(x) \middle\| \frac{p_d(x) + p_g(x)}{2} \right).$$

Podemos escribir la ecuación anterior utilizando la divergencia de Jensen-Shannon:

$$C(G) = -\log 4 + 2D_{\text{JS}}(p_d(x) \| p_g(x)).$$

Dado que D_{JS} es no negativa y vale cero en el caso de que las distribuciones sean iguales, tenemos que $-\log 4$ es el mínimo global de $C(G)$ y la única solución es $p_g = p_d$. \square

3.2.2. Convergencia

Una vez hemos probado que el óptimo global del juego minimax existe y es único, necesitamos una manera de llegar a él, es decir, demostrar que existe un algoritmo que converge a dicho óptimo.

Teorema 3.3. *Si G y D tienen suficiente capacidad y en cada paso del Algoritmo 2 se permite al discriminador alcanzar su óptimo para dado el generador G , actualizando p_g tal que se mejore el criterio:*

$$\mathbb{E}_{x \sim p_d} [\log(D_G^*(x))] + \mathbb{E}_{x \sim p_g} [\log(1 - D_G^*(x))],$$

entonces p_g converge a p_d .

Demostración. Sea $V(G, D) = U(p_g, D)$, dependiente de p_g , con $U(p_g, D)$ convexa. Las subderivadas del supremo de una función convexa incluyen la derivada de la función donde se alcanza el máximo. Es decir, si $f(x) = \sup_{a \in \mathcal{A}} f_\alpha(x)$ y $f_\alpha(x)$ es convexa para todo α , entonces:

$$\partial f_\beta(x) \in \partial f,$$

si $\beta = \arg \sup_{\alpha \in \mathcal{A}} f_\alpha(x)$. Esto es equivalente a actualizar p_g según el descenso en la dirección del gradiente para el discriminador óptimo dado el generador G . Como hemos visto en el Teorema 3.2 $\sup U(p_g, D)$ es convexo en p_g y con óptimo global único, por tanto, con actualizaciones lo suficientemente pequeñas de p_g , obtendremos la convergencia deseada. \square

Es importante resaltar que, aunque el juego minimax tiene estas excepcionales propiedades teóricas, en la práctica es necesario hacer una ligera modificación. Al principio del entrenamiento, cuando G aún tiene bajas prestaciones, D rechaza sus muestras con mucha confianza, lo que provoca que $\log(1 - D(G(z)))$ se sature con facilidad. Por ello, es frecuente entrenar G para que maximice $\log D(G(z))$, que converge al mismo equilibrio pero permite entrenar con mayor estabilidad.

Finalmente, un problema común con el que nos encontramos al entrenar redes generativas antagónicas, es que es muy complicado comparar los resultados obtenidos. En general, no hay una manera infalible de comprobar que las imágenes generadas por una red son mejores que las generadas por

otra. Se han propuesto varias medidas que tratan de cuantificar de manera objetiva la calidad de las muestras generadas, pero por el momento es un tema de investigación activo. Un buen artículo que inspecciona el estado del arte es [14].

Capítulo 4

Generación de arte

En esta sección nos vamos a apoyar sobre el artículo [16] para implementar en Python, haciendo uso de PyTorch, una red convolucional generativa antagónica profunda (**DCGAN** o Deep Convolutional Generative Adversarial Network en inglés), con la cual vamos a generar imágenes de cuadros realistas.

4.1. DCGAN

Las DCGANs son una de las extensiones más exitosas de las GANs tradicionales. Una de sus características principales es el uso de redes convolucionales con arquitecturas relativamente sencillas para el discriminador y redes con convoluciones fraccionales para el generador.

Una parte fundamental en el generador es obtener una imagen a partir de un vector de ruido aleatorio. Para ello, es necesario ir aumentando la dimensión progresivamente. Veamos pues, en primer lugar, una manera sencilla de alcanzar este objetivo. Podríamos simplemente duplicar las filas y columnas según un cierto factor. En la Figura 4.1 mostramos un ejemplo de esta técnica con factor 2.

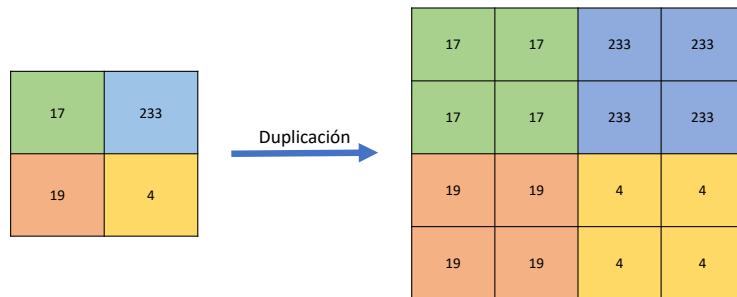


Figura 4.1: Aumento del tamaño de una imagen mediante duplicación.

También podríamos conseguir aumentar el tamaño de entrada mediante técnicas algo más sofisticadas a la descrita anteriormente, como por ejemplo la interpolación bilineal o bicúbica. Dicho esto, la técnica que mejores resultados ha cosechado es la convolución fraccional, que permite que los parámetros que llevan a cabo las transformaciones sean aprendidos por la propia red.

Es importante no confundir las convoluciones fraccionales o traspuestas con las deconvoluciones, que tienen como objetivo recuperar la entrada de una convolución, es decir, invertir la operación de convolución. Una convolución fraccional tan solo busca recuperar las dimensiones originales. Para ello, realiza la convolución usual pero añadiendo ceros de manera ingeniosa. En la Figura 4.2 describimos visualmente la manera de añadir dicho relleno de ceros. Es importante notar que la cuadrícula se mueve siempre como una convolución con stride 1. Podemos observar que el stride s de la convolución fraccional equivale a un stride de $\frac{1}{s}$ en la convolución usual y que un valor de padding de 0 equivale a un padding de $\frac{K+1}{2}$.

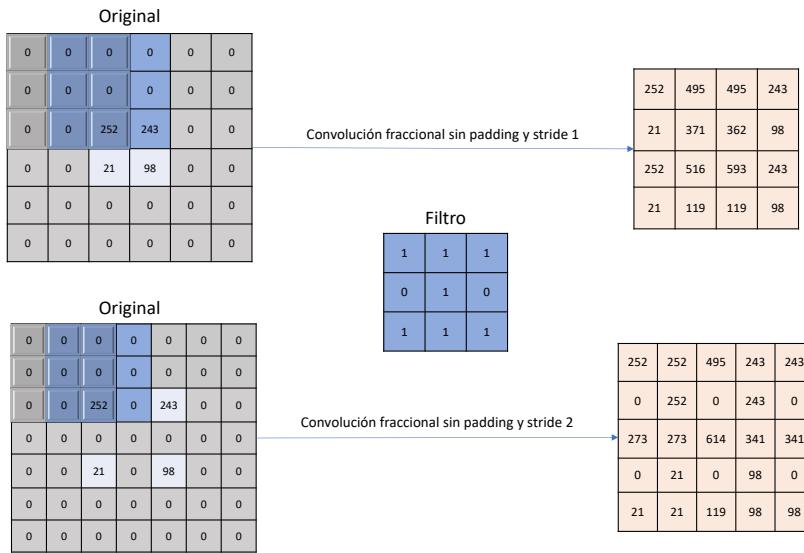


Figura 4.2: Aumento del tamaño de una imagen mediante convolución fraccional.

Fijémonos ahora en el artículo [16], donde los autores proporcionan una serie de recomendaciones empíricas para incrementar la estabilidad y obtener imágenes de mayor calidad. Los puntos más relevantes que se pueden extraer son:

- En el discriminador, utilizar convoluciones con stride en lugar de convoluciones con pooling.
- En el generador, utilizar convoluciones traspuestas con stride en lugar de convoluciones con pooling.
- No utilizar capas totalmente conectadas.
- Utilizar funciones de activación ReLU en todas las capas del generador salvo en la última, en la que utilizar tanh.

- Utilizar funciones de activación LeakyReLU en todas las capas del discriminador.
- Utilizar batch normalization tanto para el generador como el discriminador.
- Inicializar los pesos de ambas redes según una distribución normal.

Nuestro objetivo es la generación de cuadros realistas y, para ello, nos hace falta una base de datos de gran tamaño con cuadros de distintos artistas, épocas y estilos. En una competición de Kaggle¹ hemos encontrado un conjunto de datos con más de 100000 imágenes, ocupando aproximadamente 49 GB en disco. Una muestra de dichas imágenes se puede observar en la Figura 4.3.



Figura 4.3: Imágenes del conjunto de datos original sin preprocesar.

4.1.1. Preprocesado

Como paso previo a la definición de la arquitectura de las redes y su posterior entrenamiento, ha sido necesario procesar las imágenes, paso que nos ha servido también para formarnos una idea de cómo es el conjunto de datos.

En un principio teníamos imágenes etiquetadas con el nombre del autor y estilo de cada cuadro. Dado que dicha información no es relevante para nuestro propósito, la hemos descartado (se podría haber aprovechado, pero excede el alcance de este trabajo). Por otro lado, teníamos las imágenes divididas en conjuntos de entrenamiento y test, que hemos unificado, ya que nuestro problema pertenece al ámbito del aprendizaje no supervisado. Finalmente nos hemos encontrado con imágenes corruptas que hemos tenido que eliminar, por ejemplo, imágenes descargadas incorrectamente, de tamaños excesivamente grandes o en formatos incorrectos.

Una vez limpio el conjunto de datos, hemos realizado algunas transformaciones para que, posteriormente, nuestras redes reciban elementos de entrada uniformes. En primer lugar hemos escalado los tamaños y las proporciones,

¹<https://www.kaggle.com/c/painter-by-numbers/data>.

pasando así de cuadros de todos los tamaños y de distintas formas a cuadrados de 64×64 píxeles. Posteriormente, realizamos un recorte centrado, normalizamos y finalmente las cargamos como tensores, que es el formato utilizado por PyTorch. En la Figura 4.4 se puede observar una muestra de las imágenes preprocesadas.

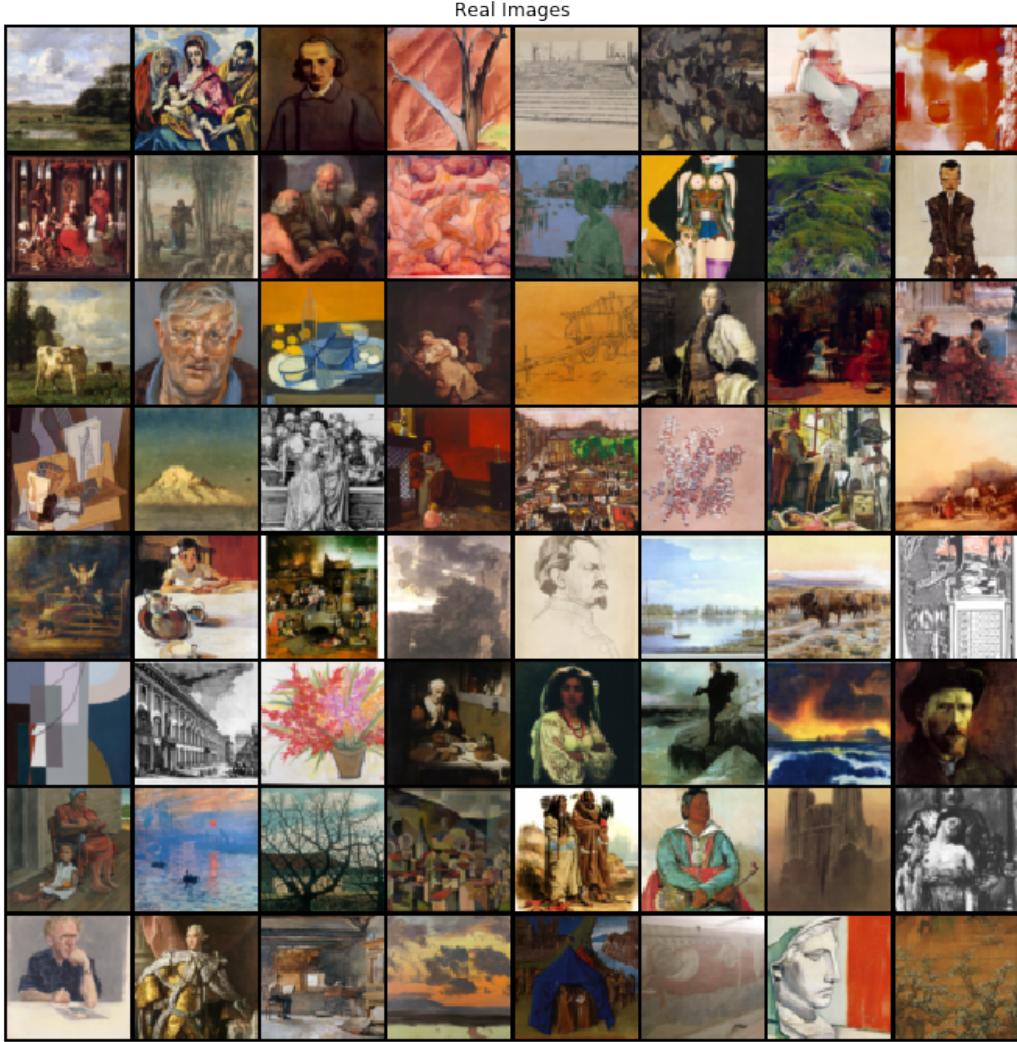


Figura 4.4: Imágenes del conjunto de datos original con el preprocesado realizado.

4.1.2. Arquitectura

Basándonos en las recomendaciones mencionadas anteriormente, hemos decidido utilizar la arquitectura que mostramos en la Figura 4.5, con una red convolucional sencilla para el discriminador y una red de convoluciones fraccionales para el generador. Además de las recomendaciones mencionadas anteriormente, hemos utilizado una técnica conocida como suavizado de etiquetas por un lado (one sided label smoothing en inglés [17]). Consiste en modificar el valor de la etiqueta correspondiente a las imágenes reales de 1 a 0,9 para reducir la confianza de la red discriminadora al hacer las predicciones, ganando así estabilidad. Para consultar los detalles a más bajo nivel,

recomendamos visitar la página correspondiente al código en el repositorio creado para este trabajo².

Hemos inicializado los pesos de ambas redes según una distribución normal $\mathcal{N}(0, 0,0002)$ y seleccionado un tamaño de batch de 128. El algoritmo de optimización por el que hemos optado es Adam [12], seleccionando una tasa de aprendizaje de 0,0002 y $\beta = (0,5, 0,999)$. Dado que el entrenamiento es muy costoso, hemos considerado prudente entrenar durante 30 épocas como máximo, guardando no obstante los modelos en cada época para después poder volver a ellos y realizar comparaciones en la calidad de las imágenes generadas. Como ya hemos mencionado, medir la calidad de las imágenes generadas es un tema delicado, por lo que las comparaciones que hemos realizado han sido a nivel visual. Para poder hacer las comparaciones algo más fiables, hemos fijado un vector de ruido aleatorio al principio y generado una muestra a partir de él cada 200 batches y al final de cada época, obteniendo de este modo una visualización de la evolución del entrenamiento³.

4.1.3. Recursos

Como hemos ido comentando a lo largo de todo el trabajo, el entrenamiento de las GANs no es para nada sencillo, siendo necesaria una cuidadosa elección del conjunto de datos, de la arquitectura de la red y de los hiperparámetros. Sin embargo esto es tan solo una parte de la dificultad. También es necesario disponer de unas capacidades de cómputo generosas. Ha sido imprescindible el uso de un ordenador con GPU compatible con CUDA y una gran cantidad de memoria RAM. En concreto el equipo que hemos utilizado contaba con una GPU Nvidia Quadro P5000 y 32GB de RAM. Con esta configuración, las 30 épocas de entrenamiento se prolongaron durante aproximadamente 24 horas. Para contrastar, también hemos realizado pruebas entrenando con CPU en un portátil de modestas prestaciones, observando un rendimiento aproximadamente 20 veces menor.

4.2. Resultados

En la Figura 4.6 mostramos medio batch de imágenes generadas con nuestra GAN a partir de un vector de entrada aleatorio, sin realizar ningún tipo de selección ni procesado. Observamos que, en general y a primera vista, algunas podrían hacerse pasar por obras de arte, incluso podríamos distinguir paisajes, retratos o composiciones abstractas. Sin embargo, tras una inspección más cuidadosa, es claro que las imágenes que no pueden ser catalogadas como abstractas carecen de detalles o presentan formas inusuales, bastante típicas en las imágenes generadas con GANs. Con todo, consideramos que las imágenes obtenidas son visualmente agradables.

²https://github.com/ant-mak/tfm/blob/master/src/tfm_teci_antonmakarov_gan-paintings.ipynb.

³Una animación de la evolución de la muestra desde ruido aleatorio hasta imágenes realistas se puede ver en el repositorio.

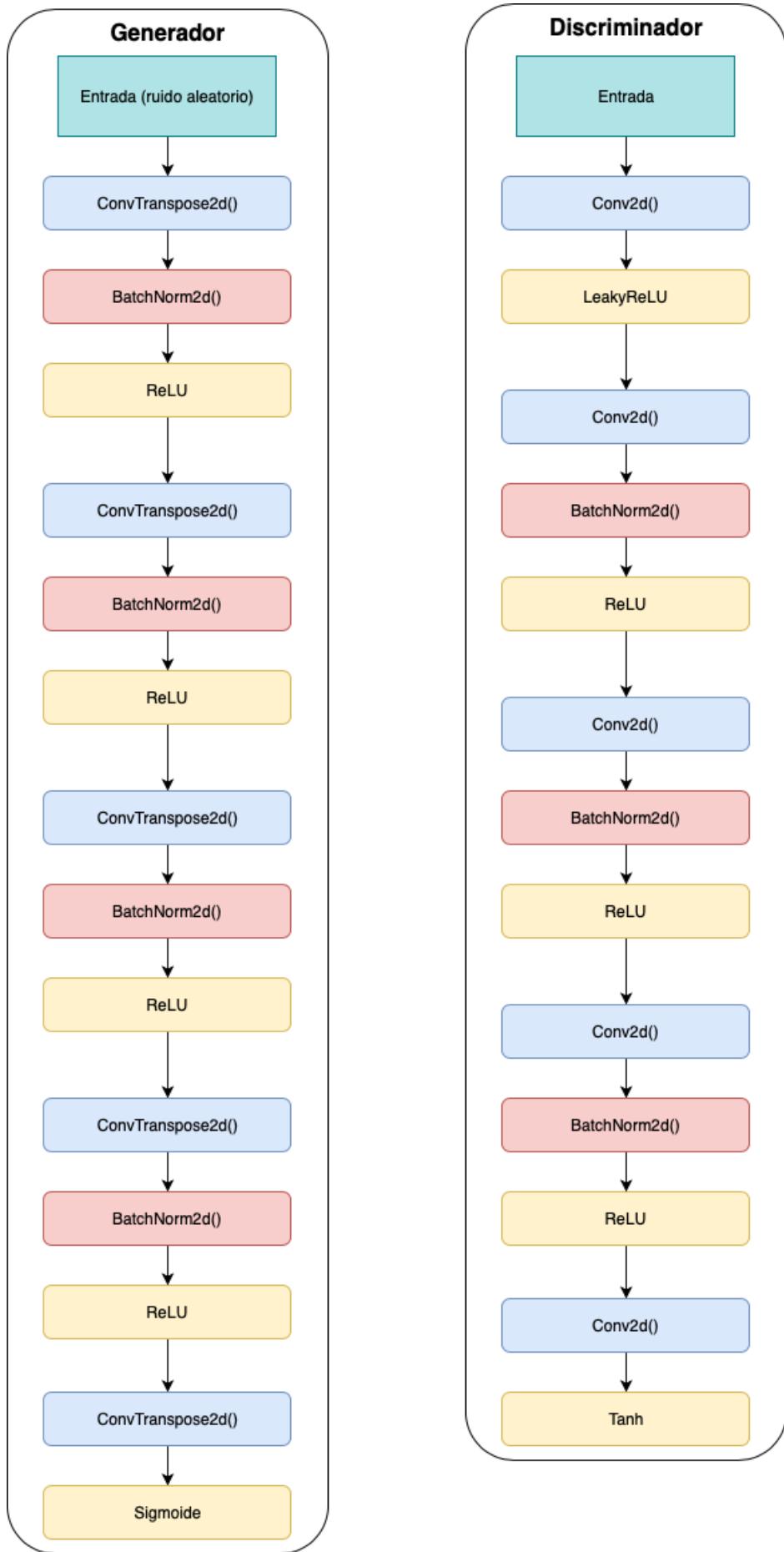


Figura 4.5: Arquitectura de las redes (Cambiar, poner tamaños de entrada y salida).

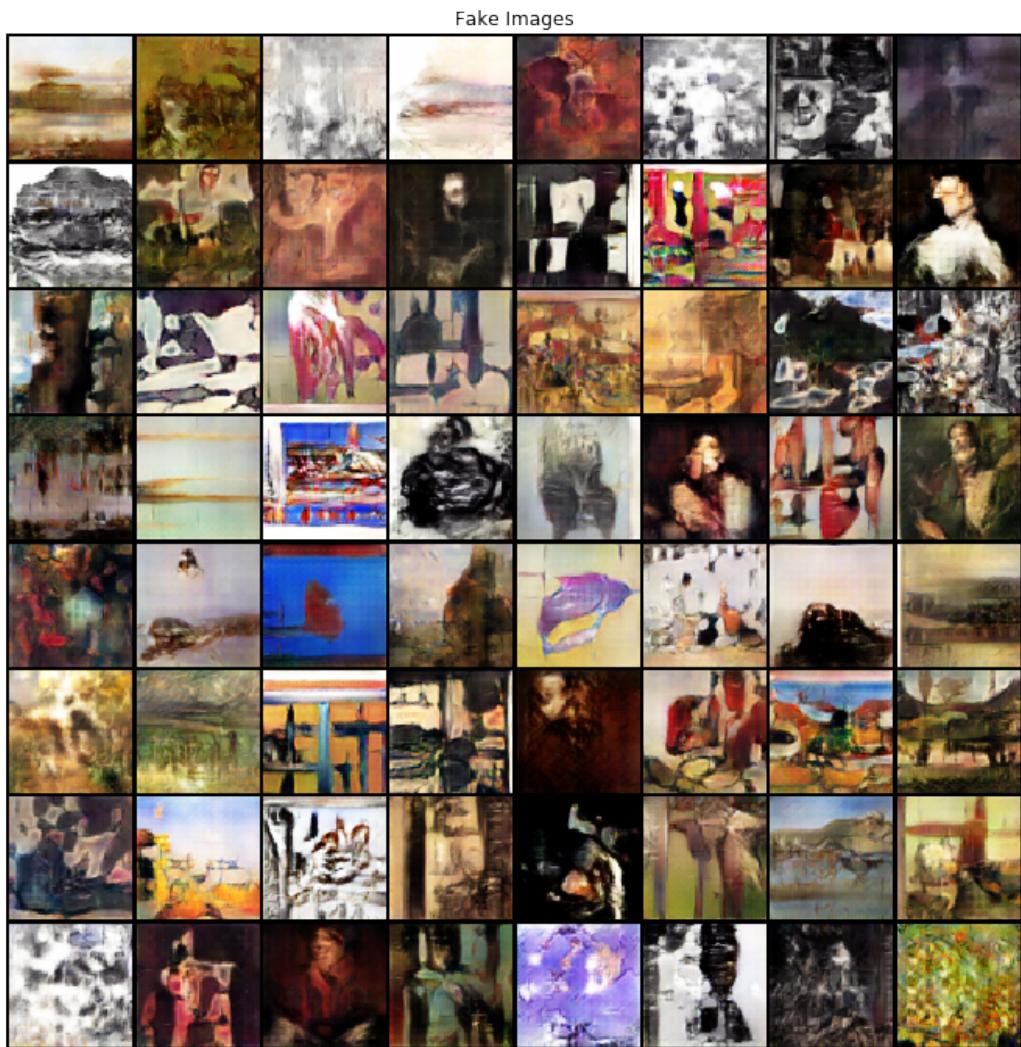


Figura 4.6: Imágenes generadas mediante la DCGAN implementada después de 30 épocas.

Capítulo 5

Conclusión

Las redes generativas antagónicas son una herramienta de gran interés tanto desde el punto de vista teórico como práctico. Por un lado introducen una nueva manera de pensar en el mundo del aprendizaje automático, generando una cantidad de investigación casi sin precedentes. Por otro, los modelos generativos pueden ser muy útiles para la industria, pudiendo ser empleados en, por ejemplo, la separación de contexto y estilo en imágenes [11], el entrenamiento de algoritmos igualitarios [19] o la generación de texto [18].

Nuestro objetivo con este trabajo ha sido la construcción de una primera aproximación al campo de las GANs, incluyendo una breve introducción del contexto del problema, el desarrollo de la idea conceptual, una reseña de la teoría involucrada y, por último, una implementación práctica, que ha sido fundamental para consolidar los conceptos más intrincados, aprender a trabajar con PyTorch y utilizar una GPU para llevar a cabo el entrenamiento.

Los resultados que hemos obtenido en este trabajo tienen aún mucho potencial de mejora y nos gustaría seguir trabajando en ello. Algunas de las opciones que barajamos son la inclusión de etiquetas sobre el estilo de los cuadros en el modelo, que ya tenemos en nuestro conjunto de datos pero no hemos utilizado. Con esto, esperamos poder generar imágenes restringidas al estilo que queramos, siguiendo una metodología similar a [15]. También sería interesante tratar de aumentar la resolución, ya que actualmente tenemos imágenes de 64×64 . Para ello podríamos tratar de incorporar el entrenamiento progresivo descrito en [10]. Por último, sería interesante seleccionar los hiperparámetros principales de los modelos de una manera más metódica, utilizando alguno de los métodos descritos en [3].

También nos gustaría investigar en mayor profundidad el artículo [1], cuyo contenido matemático es excelente y [2], que aporta gran cantidad de explicaciones conceptuales e intuitivas. Así mismo, queremos continuar con el estudio de técnicas y arquitecturas más avanzadas, por ejemplo [11, 18], que están a la vanguardia en el tema de transferencia de estilo y generación de texto respectivamente. Finalmente, no queremos dejar sin mencionar el campo de los ejemplos antagónicos [8], un tema estrechamente relacionado con las GANs, que introduce el concepto de ataque a algoritmos de aprendizaje automático. Por ejemplo, se consigue que un algoritmo clasifique una imagen de un perro como taza de café añadiendo tan solo un poco de ruido imperceptible al ojo humano.

Bibliografía

- [1] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [2] D. Bau, J.-Y. Zhu, H. Strobelt, B. Zhou, J. B. Tenenbaum, W. T. Freeman, and A. Torralba. Gan dissection: Visualizing and understanding generative adversarial networks. *arXiv preprint arXiv:1811.10597*, 2018.
- [3] J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pages 2546–2554, 2011.
- [4] K. Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, (36):193–302, 1980.
- [5] I. Goodfellow. Nips 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016.
- [6] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [8] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [9] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [10] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- [11] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019.
- [12] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

- [13] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [14] M. Lucic, K. Kurach, M. Michalski, S. Gelly, and O. Bousquet. Are gans created equal? a large-scale study. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 700–709. Curran Associates, Inc., 2018.
- [15] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [16] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [17] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *Advances in neural information processing systems*, pages 2234–2242, 2016.
- [18] L. Yu, W. Zhang, J. Wang, and Y. Yu. Seqgan: Sequence generative adversarial nets with policy gradient. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [19] B. H. Zhang, B. Lemoine, and M. Mitchell. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 335–340. ACM, 2018.
- [20] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.