

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

Resumen

Abstract

Abstract

Índice general

1. Introducción	7
2. Preliminares	8
2.1. Conceptos básicos	8
2.1.1. Aprendizaje supervisado	8
2.1.2. Aprendizaje no supervisado	9
2.1.3. Aprendizaje por refuerzo	9
2.2. Redes neuronales	9
2.3. Aprendizaje profundo	10
2.3.1. Redes neuronales convolucionales	10
3. Redes generativas antagónicas	15
3.1. Idea general	15
3.2. Bases teóricas	16
3.3. Otros modelos generativos	17
4. Generación de arte	18
4.1. DCGAN	18
4.1.1. Preprocesado	19
4.1.2. Arquitectura	19
4.1.3. Recursos	20
4.2. Resultados	20
5. Conclusión	23
Bibliografía	24

Capítulo 1

Introducción

Introducción¹ [4], [2], [7], [3]

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

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. Conceptos básicos

El aprendizaje automático tiene como objetivo principal el desarrollo de técnicas que permitan aprender a las máquinas 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. 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.

2.3. 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.

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.

2.3.1. 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 [1] que, posteriormente, sería tomado por Yann LeCun, 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 2.1. *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.*

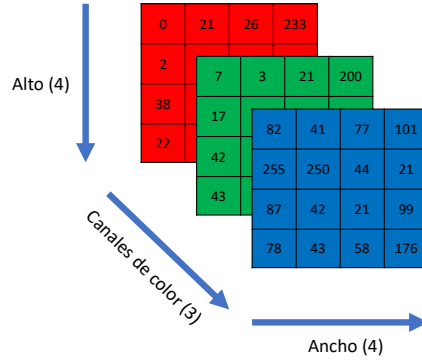


Figura 2.1: Representación esquemática de una imagen.

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 predeterminado, 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. 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.

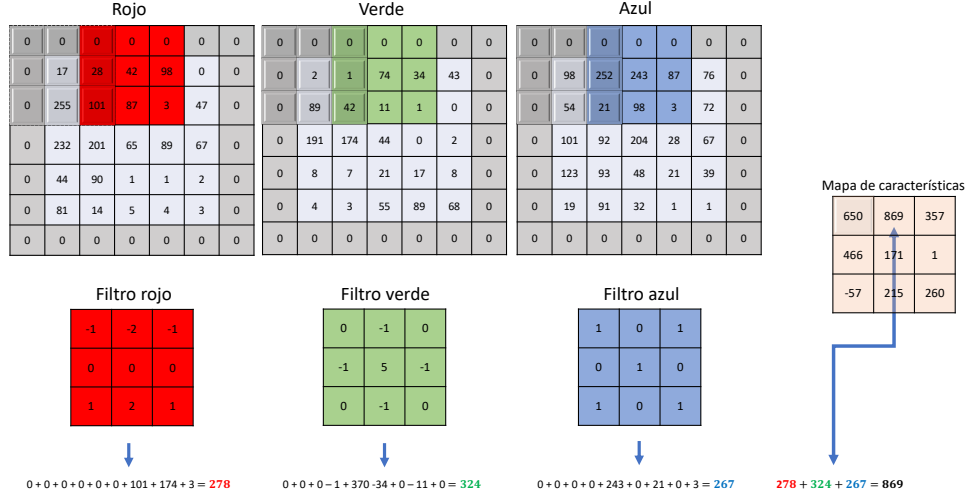


Figura 2.2: Diagrama ilustrativo de la operación de convolución.

Observación 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.

Recientemente se ha incorporado una técnica para mejorar la estabilidad en el entrenamiento de redes convolucionales, llamada **Batch Normalization** [5]. 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 [5].

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 excesiva-

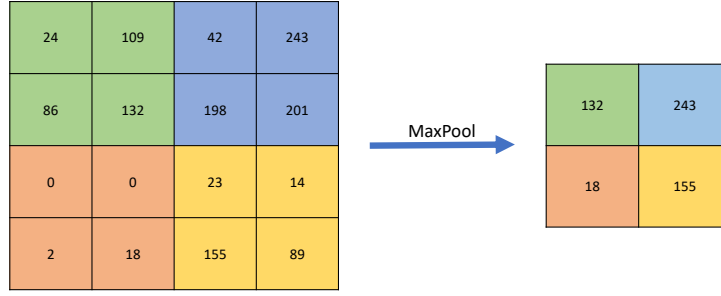


Figura 2.3: Diagrama ilustrativo de reducción de muestreo mediante máximo.

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: $y_i = \gamma x_i + \beta$;
-

mente 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 2.2. 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*:

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

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 contruidos los mapas de características, se aplanan 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.

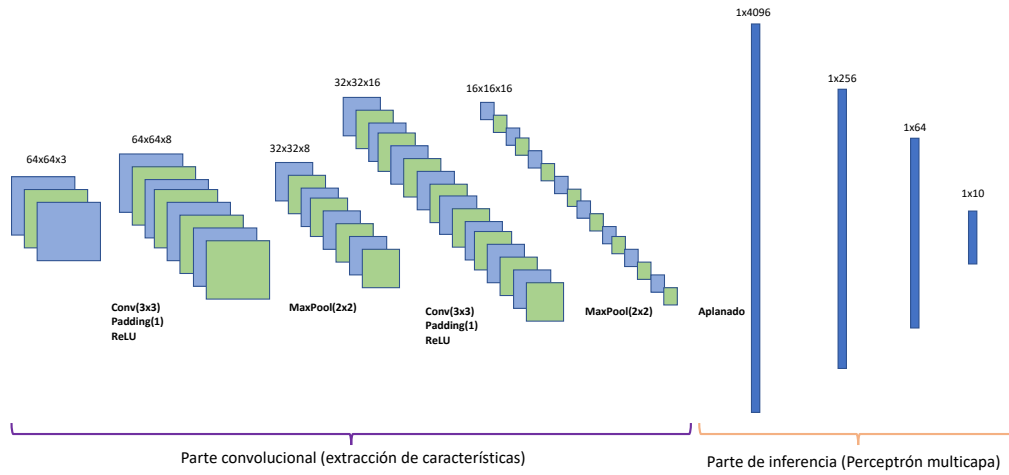


Figura 2.4: Arquitectura básica de una red convolucional

Capítulo 3

Redes generativas antagónicas

En este capítulo describiremos las ideas principales, tanto teóricas como prácticas de las redes generativas antagónicas, generative adversarial networks o GANs en inglés. Fueron propuestas por primera vez por Ian Goodfellow en 2014 [4], mezclando conceptos de aprendizaje automático no supervisado, supervisado y teoría de juegos. Desde entonces han suscitado una actividad investigadora muy importante, con aplicaciones en prácticamente todos los campos relacionados con el aprendizaje automático.

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 obtener imágenes a partir de ruido aleatorio con distribución $p_g(z)$, 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 sepa discernir los billetes falsos de los verdaderos, obteniendo así los ladrones una falsificación realista. Una representación esquemática se puede ver en la Figura 3.1. Podemos traducir esta idea a términos matemáticos

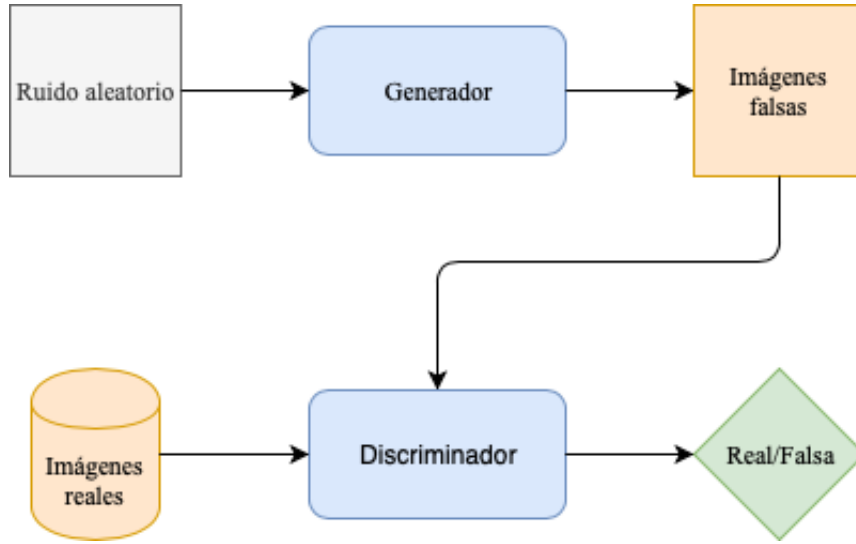


Figura 3.1: Diagrama conceptual de las redes generativas antagónicas.

de la siguiente manera:

Para aprender la distribución p_g del generador sobre los datos x , definimos una distribución sobre las variables de ruido de entrada $p_z(z)$. Mediante un perceptrón multicapa con parámetros θ_g definimos la función $G(z; \theta_g)$. Definimos también mediante otro perceptrón multicapa $D(x; \theta_d)$ la probabilidad de que x es una muestra de los datos y no del generador. Entrenamos D para que maximice la probabilidad de asignar la etiqueta correcta tanto a los ejemplos de entrenamiento como a los generados. Al mismo tiempo, entrenamos G para que minimice $\log(1 - D(G(z)))$. Es decir, tenemos el siguiente juego minimax:

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

En el Algoritmo 2 describimos lo anterior de manera más concisa y damos una primera idea de cómo se lleva a cabo el entrenamiento.

3.2. Bases teóricas

En esta sección realizaremos un análisis de la teoría que hay detrás de las redes generativas antagónicas. Mostraremos que el criterio de entrenamiento nos permite recuperar la distribución de los datos si damos a G y D capacidad suficiente.

Teorema 1. *Para un generador G fijo, el discriminador D óptimo es:*

$$D_G^*(x) = \frac{p_{\text{datos}}(x)}{p_{\text{datos}}(x) + p_g(x)} \quad (3.2)$$

Demostración. El discriminador busca maximizar su función de utilidad, dada por $V(G, D)$

$$\begin{aligned} V(G, D) &= \int_x p_{\text{datos}}(x) \log(D(x)) dx + \int_z p_z(z) \log(1 - D(g(z))) dz \\ &= \int_x p_{\text{datos}}(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx \end{aligned} \quad (3.3)$$

Algoritmo 2: Entrenamiento de una red generativa antagónica genérica.

```
1 for Iteraciones de entrenamiento do
2   Tomar muestra  $(z^{(1)}, z^{(2)}, \dots, z^{(m)})$  de tamaño  $m$  de  $p_g(z)$ ;
3   Tomar muestra  $(x^{(1)}, x^{(2)}, \dots, x^{(m)})$  de tamaño  $m$  de  $p_d(x)$ ;
4   Actualizar el discriminador ascendiendo su gradiente:

      
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(x^i) + \log (1 - D(G(z^i)))]$$


5   Actualizar el generador ascendiendo su gradiente:

      
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^i)))$$


6 end
```

□

3.3. Otros modelos generativos

Variational autoencoders, noise contrastive estimation etc.

Capítulo 4

Generación de arte

En esta sección nos vamos a apoyar sobre el artículo [7] para implementar en `Python`, haciendo uso de `PyTorch` una red convolucional profunda generativa antagónica, con la cual vamos a tratar de generar imágenes de cuadros realistas.

4.1. DCGAN

Las redes convolucionales profundas generativas antagónicas (deep convolutional generative adversarial networks, DCGAN) son una extensión de las GANs tradicionales, que implementan una serie de recomendaciones para incrementar la estabilidad y obtener imágenes de mejor calidad. Los puntos más relevantes que se pueden extraer de [7] son:

- En el discriminador, utilizar convoluciones con stride en lugar de capas de pooling.
- En el generador, utilizar convoluciones fraccionales con stride en lugar de capas de pooling.
- Utilizar BatchNorm tanto para generador como discriminador.
- No utilizar capas fully connected.
- Utilizar funciones de activación ReLU en el generador, salvo para la última capa, en la que se propone usar tanh.
- Utilizar funciones de activación LeakyReLU en el discriminador.
- Inicializar los pesos de ambas redes con una distribución normal.

RECORDAR QUE STRIDE Y POOL SON PARECIDOS Y DESCRIBIR LEAKYRELU.

El objetivo que nos hemos propuesto es la generación de cuadros realistas, para ello, es necesario un amplio conjunto de datos con cuadros de distintos artistas, épocas y estilos. Hemos encontrado en una competición de Kaggle un conjunto de datos con las características deseadas, con más de 100000 imágenes, ocupando aproximadamente 49 GB. Una muestra de dichas imágenes se puede observar en la Figura 4.1

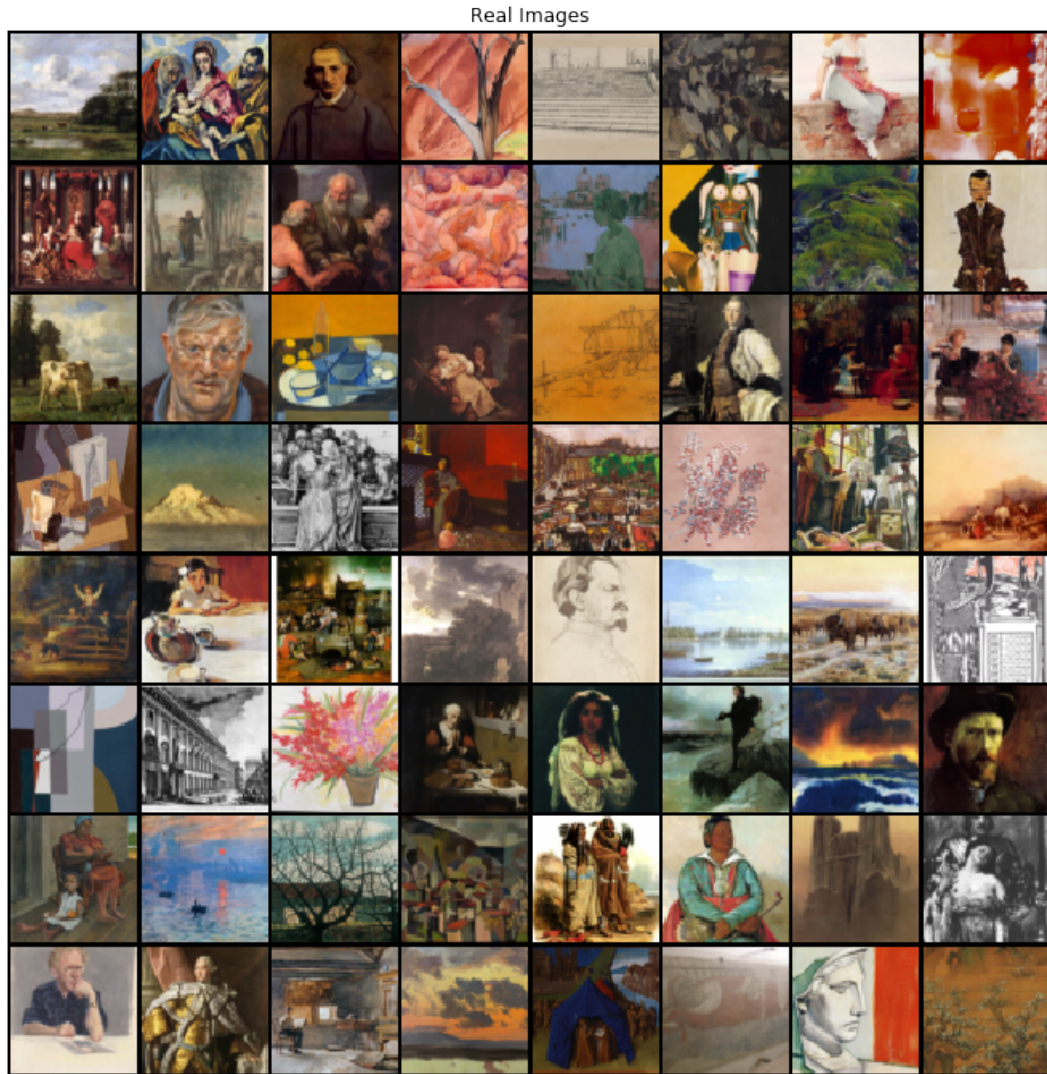


Figura 4.1: Imágenes del conjunto de datos.

4.1.1. Preprocesado

El conjunto de datos tenía algunos defectos, como por ejemplo imágenes descargadas incorrectamente o imágenes corruptas en el propio conjunto de datos. Se ha tomado la decisión de eliminar dichas imágenes. ya que son pocas y no hay una manera directa de recuperarlas.

Una vez limpio el conjunto de datos, es necesario realizar algunas transformaciones para que, posteriormente, nuestras redes neuronales puedan utilizar las imágenes. En primer lugar hemos homogeneizado los tamaños y las proporciones. Pasamos así de cuadros de todos los tamaños y de distintas formas a cuadrados de 64×64 píxeles. Posteriormente, realizamos un recorte centrado y finalmente las cargamos como tensores.

4.1.2. Arquitectura

Siguiendo las recomendaciones del artículo [7], hemos decidido utilizar la arquitectura que se muestra en la Figura 4.1.2. También hemos inicializado

los pesos con una distribución normal $\mathcal{N}()$ y hemos seleccionado un tamaño de batch de 128. La red se ha entrenado durante 30 épocas, utilizando el algoritmo de optimización Adam [6] con una tasa de aprendizaje de 0,0002 y $\beta = (0,5, 0,999)$.

4.1.3. Recursos

Cabe destacar que ha sido imprescindible el uso de un ordenador con GPU compatible con CUDA, en concreto se ha utilizado una Nvidia Quadro P5000. El entrenamiento de principio a fin ha tardado aproximadamente 24 horas. También se han realizado pruebas con CPU en un portátil de prestaciones modestas y se ha observado que el rendimiento era unas 20 veces menor.

4.2. Resultados

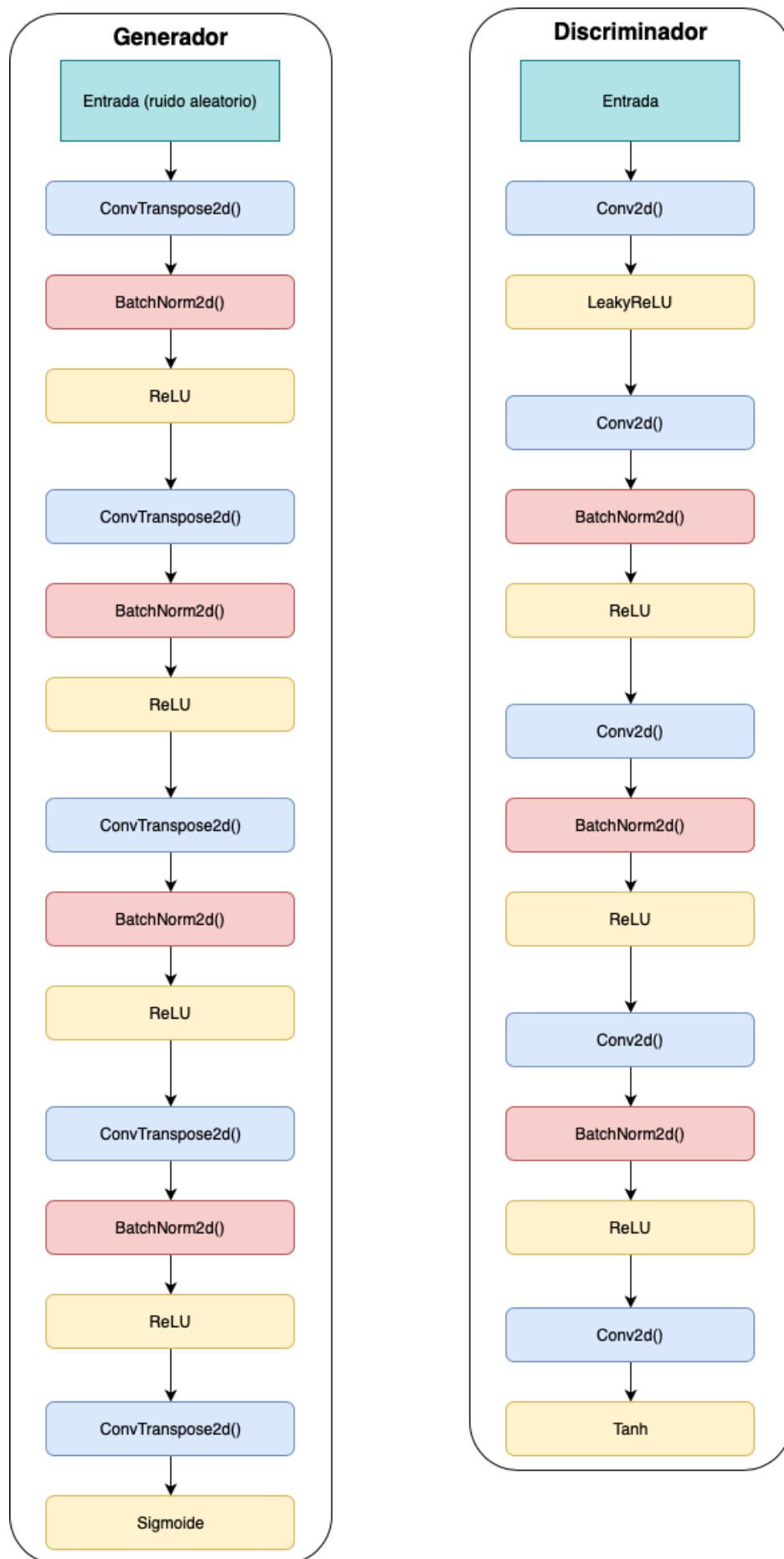


Figura 4.2: Arquitectura de las redes (Cambiar, poner tamaños de entrada y salida).

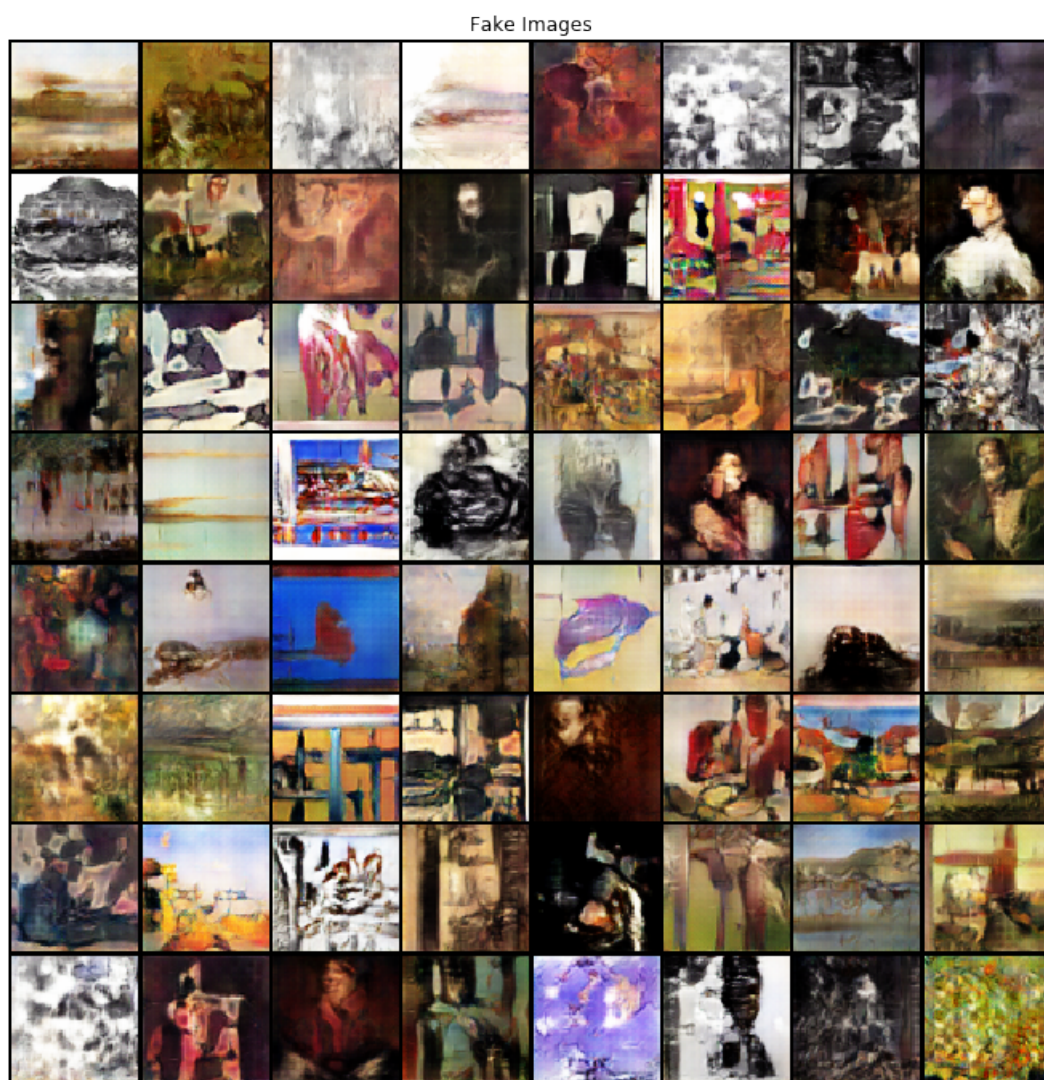


Figura 4.3: Imágenes generadas mediante la DCGAN implementada después de 30 épocas de entrenamiento.

Capítulo 5

Conclusión

Bibliografía

- [1] 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.
- [2] I. Goodfellow. Nips 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016.
- [3] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [4] 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.
- [5] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [6] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [7] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.