

# ÍNDICE

---

<b>PREFACIO .....</b>	<b>XI</b>
<b>CAPÍTULO 1. INTRODUCCIÓN.....</b>	<b>1</b>
1.1 PRELIMINARES .....	1
1.2 BREVE REVISIÓN HISTÓRICA.....	7
1.3 POTENCIALES DESTINATARIOS DEL LIBRO .....	9
1.4 HERRAMIENTAS Y CONJUNTOS DE DATOS.....	11
1.4.1 Herramientas, <i>frameworks</i> , librerías y servicios cognitivos .....	11
1.4.2 Conjuntos de datos.....	16
1.5 ORGANIZACIÓN DEL LIBRO .....	28
<b>CAPÍTULO 2. COMPUTACIÓN NUMÉRICA.....</b>	<b>31</b>
2.1 INTRODUCCIÓN .....	31
2.2 TENSORES .....	31
2.3 DESBORDAMIENTO Y SUBDESBORDAMIENTO .....	32
2.4 MÉTODOS DE OPTIMIZACIÓN .....	33
2.4.1 Gradiente Descendente Estocástico.....	37
2.4.2 Propagación de la Raíz Media Cuadrática.....	42
2.4.3 Estimación del Momento Adaptativo.....	43
2.4.4 Otros métodos de optimización .....	43
2.5 FUNCIONES DE ACTIVACIÓN NO LINEALES Y UNIDADES LINEALES.....	44
2.6 FUNCIONES DE PÉRDIDA .....	50
2.6.1 Clasificación .....	51
2.6.2 Regresión .....	53
<b>CAPÍTULO 3. REDES NEURONALES PROFUNDAS.....</b>	<b>55</b>
3.1 INTRODUCCIÓN .....	55
3.2 FUNDAMENTOS GENERALES .....	55
3.3 EL PERCEPTRÓN.....	58
3.3.1 Arquitectura del perceptrón.....	58
3.3.2 Entrenamiento del perceptrón.....	61
3.3.3 El perceptrón para problemas multiclase.....	64
3.4 LA RED DE RETROPROPAGACION .....	66
3.4.1 Arquitectura de la red.....	66
3.4.2 Entrenamiento por retropropagación .....	69

3.5 REDES DE CREENCIA/BAYESIANAS PROFUNDAS .....	76
<b>CAPÍTULO 4. OPERACIONES REDES NEURONALES CONVOLUCIONALES I .....</b>	<b>79</b>
4.1 INTRODUCCIÓN .....	79
4.2 OPERACIÓN DE CONVOLUCIÓN .....	80
4.2.1 Sin relleno con ceros y una unidad de desplazamiento .....	87
4.2.2 Relleno con ceros y una unidad de desplazamiento .....	88
4.2.3 Relleno para obtener la misma dimensión de salida .....	89
4.2.4 Relleno completo para obtener una mayor dimensión de salida.....	90
4.2.5 Sin relleno con ceros y desplazamientos superiores a la unidad .....	91
4.2.6 Relleno con ceros y desplazamientos superiores a la unidad .....	91
4.3 AGRUPAMIENTO ( <i>POOLING</i> ) .....	93
4.4 CONVOLUCIÓN ARITMÉTICA TRANSPUESTA .....	98
4.4.1 Transposición no <i>zero-padding</i> y una unidad de desplazamiento .....	101
4.4.2 Transposición <i>zero-padding</i> y una unidad de desplazamiento .....	102
4.4.3 Transposición <i>half (same) padding</i> .....	103
4.4.4 Transposición <i>full padding</i> .....	103
4.4.5 Transposición no <i>zero-padding</i> y desplazamientos superiores a la unidad .....	104
4.4.6 Transposición <i>zero-padding</i> y desplazamientos superiores a la unidad .....	105
<b>CAPÍTULO 5. OPERACIONES REDES NEURONALES CONVOLUCIONALES II .....</b>	<b>109</b>
5.1 INTRODUCCIÓN .....	109
5.2 CONVOLUCIONES DILATADAS .....	109
5.3 CONVOLUCIÓN NO LINEAL .....	112
5.4 CONVOLUCIONES 1D, 2D y 3D .....	115
5.5 SOBREAJUSTE, <i>WEIGHT DECAY</i> Y <i>DROPOUT</i> .....	123
5.6 NORMALIZACIÓN .....	129
5.7 <i>SOFTMAX</i> .....	133
<b>CAPÍTULO 6. MOTIVACIÓN DE LAS REDES NEURONALES CONVOLUCIONALES .....</b>	<b>135</b>
6.1 INTRODUCCIÓN .....	135
6.2 COMPARTICIÓN DE PARÁMETROS .....	135
6.3 BASE NEUROCIÉNTIFICA DE LAS CNN .....	140
<b>CAPÍTULO 7. ARQUITECTURAS DE LAS REDES NEURONALES CONVOLUCIONALES I .....</b>	<b>149</b>
7.1 INTRODUCCIÓN .....	149
7.2 ORGANIZACIÓN DE CAPAS EN LAS CNN .....	149
7.3 CAPAS <i>INCEPTION</i> .....	150
7.4 ALEXNET .....	156
7.4.1 Niveles de activación y visualización en las capas de la red .....	158
7.4.2 Agrupamiento piramidal espacial .....	160
7.5 VGGNET .....	161
7.6 LENET .....	163
7.7 RESNET .....	164
<b>CAPÍTULO 8. ARQUITECTURAS DE LAS REDES NEURONALES CONVOLUCIONALES II .....</b>	<b>173</b>
8.1 INTRODUCCIÓN .....	173
8.2 GOOGLNET .....	173
8.3 OTRAS REDES CON MÓDULOS INCEPTION .....	184
8.3.1 Red Inception-V4 .....	185

8.3.2 Redes híbridas Inception-ResNet.....	187
8.4 ARQUITECTURA XCEPTION .....	191
8.5 SQUEEZENET .....	191
8.6 MÓDULOS DE CONVOLUCIONES DILATADAS PARALELAS .....	199
8.7 REDES DENSAS .....	200
8.8 REDES DENTRO DE REDES .....	202
<b>CAPÍTULO 9. ARQUITECTURAS DE LAS REDES NEURONALES CONVOLUCIONALES III .....</b>	<b>205</b>
9.1 INTRODUCCIÓN .....	205
9.2 AUTOCODIFICADORES .....	206
9.2.1 Arquitectura de un autocodificador .....	207
9.2.2 Funciones de activación y pérdida.....	213
9.2.3 Aplicaciones de los autocodificadores.....	215
9.3 REDES SIAMESAS .....	219
9.3.1 Arquitectura de las redes siamesas .....	220
9.3.2 Entrenamiento de las redes siamesas .....	223
9.3.3 Aplicaciones de las redes siamesas .....	225
9.4 REDES NEURONALES DE CÁPSULAS (CAPSNET) .....	228
<b>CAPÍTULO 10. SEGMENTACIÓN SEMÁNTICA DE IMÁGENES CON CNN .....</b>	<b>239</b>
10.1 INTRODUCCIÓN .....	239
10.2 CODIFICADOR-DECODIFICADOR (AUTOENCODER) .....	240
10.3 RED CONVOLUCIONAL TOTAL (FCN) .....	241
10.4 PARSENET .....	244
10.5 U-NET .....	244
10.6 RED PIRAMIDAL DE CARACTERÍSTICAS .....	246
10.7 PSPNET .....	250
10.8 DEEPLAB .....	251
10.8.1 DeepLabv3 .....	259
10.8.2 DeepLabv3+ .....	261
10.9 ENCNET .....	263
10.9.1 Codificación de contexto .....	264
10.9.2 Enfatización del mapa de características.....	264
10.9.3 SE-Loss .....	265
10.10 MEDIDAS DE DESEMPEÑO .....	266
<b>CAPÍTULO 11. REDES CONVOLUCIONALES PARA ANÁLISIS DE VÍDEO .....</b>	<b>267</b>
11.1 INTRODUCCIÓN .....	267
11.2 MODELOS CNN CON LSTM .....	268
11.3 REDES DE DOS FLUJOS .....	268
11.3.1 Modelo-1 .....	269
11.3.2 Modelo-2 .....	281
11.4 REDES CONVOLUCIONALES 2D, 3D Y MIXTAS .....	284
11.4.1 R2D: Convoluciones 2D sobre el vídeo completo .....	285
11.4.2 R3D: Convoluciones 3D .....	286
11.4.3 MCx y rMCx: convoluciones mixtas 3D y 2D.....	287
11.4.4 R(2+1)D: (2+1)D convoluciones .....	288
11.5 REDES DE SEGMENTOS TEMPORALES .....	289

<b>CAPÍTULO 12. DETECCIÓN DE OBJETOS EN IMÁGENES I</b> .....	<b>297</b>
12.1 INTRODUCCIÓN .....	297
12.2 COEFICIENTES DE REGRESIÓN DEL RECTÁNGULO .....	299
12.3 SOLAPAMIENTO DE REGIONES Y PRECISIÓN .....	302
12.4 ANCHOR BOXES .....	304
12.5 DETECCIÓN DE OBJETOS MULTIESCALA .....	310
12.6 R-CNN .....	311
12.7 FAST R-CNN .....	313
12.7.1 Función de pérdida <i>multi-task</i> .....	316
12.7.2 Estrategia de muestreo del <i>mini-batch</i> .....	318
12.7.3 Retropropagación a través de las capas de <i>pooling</i> de la Rol .....	319
12.7.4 Hiperparámetros SGD .....	319
12.7 FASTER R-CNN .....	320
12.8.1 Generación de regiones ( <i>anchors generation</i> ) .....	322
12.8.2 Función de pérdida ( <i>loss function</i> ) .....	322
12.8.3 Entrenamiento .....	324
12.9 MASK R-CNN .....	325
12.9.1 Entrenamiento .....	326
12.9.2 Representación de la máscara .....	326
12.9.3 Alineamiento de la Rol .....	326
12.9.4 Arquitectura de red .....	330
12.10 SSD: SINGLE SHOT MULTIBOX DETECTOR .....	331
12.10.1 Bases para el entrenamiento del SSD .....	333
12.10.2 Objetivo de entrenamiento del SSD .....	334
<b>CAPÍTULO 13. DETECCIÓN DE OBJETOS EN IMÁGENES II</b> .....	<b>335</b>
13.1 INTRODUCCIÓN .....	335
13.2 YOLO .....	335
13.2.1 Primera versión (YOLOv1) .....	336
13.2.2 Segunda versión (YOLOv2) .....	340
13.2.3 Tercera versión (YOLOv3) .....	351
13.2.4 Cuarta y quinta versiones (YOLOv4, YOLOv5) .....	355
13.3 PANET .....	357
13.3.1 Incremento de ruta abajo-arriba .....	357
13.3.2 <i>Pooling</i> adaptativo de características .....	359
13.3.3 Fusión totalmente conectada .....	360
13.4 CORNERNET/CENTERNET .....	362
13.4.1 Red <i>hourglass</i> .....	363
13.4.2 Detección de esquinas .....	365
13.4.3 Agrupación de esquinas .....	366
13.4.4 <i>Corner pooling</i> .....	367
<b>CAPÍTULO 14. DETECCIÓN DE OBJETOS EN IMÁGENES III</b> .....	<b>371</b>
14.1 INTRODUCCIÓN .....	371
14.2 OVERFEAT .....	371
14.2.1 Clasificación .....	372
14.2.2 Localización .....	375
14.2.3 Detección .....	376

14.3 RETINANET.....	377
14.4 R-FCN .....	380
14.5 FCOS.....	390
14.6 EFFICIENTDET .....	393
14.7 TRANSFORMADOR ESPACIAL .....	397
14.7.1 Red de localización .....	397
14.7.2 Rejilla de muestreo parametrizada .....	398
14.7.3 Muestreador de imagen diferenciable .....	399
<b>CAPÍTULO 15. REDES PARA DISPOSITIVOS MÓVILES .....</b>	<b>403</b>
15.1 INTRODUCCIÓN .....	403
15.2 MOBILENET (V1, V2, V3).....	403
15.2.1 MobileNetV1.....	403
15.2.2 MobileNetV2.....	405
15.2.3 MobileNetV3.....	410
15.3 SHUFFLENET .....	414
<b>CAPÍTULO 16. PLATAFORMAS PARA ESPACIOS DE BÚSQUEDA EN CLASIFICACIÓN DE IMÁGENES .....</b>	<b>419</b>
16.1 INTRODUCCIÓN .....	419
16.2 NAS .....	419
16.3 NASNET.....	424
16.4 NETADAPT .....	431
<b>CAPÍTULO 17. ALGORITMO DEEPPDREAM Y REDES GENERATIVAS ANTAGÓNICAS .....</b>	<b>437</b>
17.1 INTRODUCCIÓN .....	437
17.2 DEEPPDREAM .....	438
17.3 REDES GENERATIVAS ANTAGÓNICAS.....	449
<b>CAPÍTULO 18. REDES NEURONALES RECURRENTE, RECURSIVAS Y LSTM.....</b>	<b>459</b>
18.1 INTRODUCCIÓN .....	459
18.2 CONSIDERACIONES PRELIMINARES .....	459
18.3 REDES NEURONALES RECURRENTE .....	463
18.4 REDES NEURONALES RECURSIVAS .....	469
18.5 LONG SHORT-TERM MEMORY (LSTM) .....	470
18.5.1 LSTM Bidireccional (BiLSTM) .....	475
18.5.2 Predicción utilizando LSTM .....	477
18.5.3 Clasificación utilizando LSTM .....	480
18.6 GATED RECURRENT UNITS (GRU).....	482
18.7 ENCODER-DECODER UTILIZANDO LSTM.....	485
<b>ANEXO: RETROPROPAGACIÓN .....</b>	<b>487</b>
A.1 ESQUEMA DEL MODELO DE ENTRENAMIENTO .....	487
A.2 REGLA DELTA .....	488
A.3 MECANISMO DE RETROPROPAGACIÓN .....	490
A.3.1 Ejemplo 1: una neurona de entrada y una de salida .....	492
A.3.2 Ejemplo 2: dos neuronas de entrada y una de salida .....	493
A.3.3 Ejemplo 3: dos neuronas de entrada y dos de salida .....	496
A.4 ENTROPÍA CRUZADA Y SU DERIVADA.....	500

<b>BIBLIOGRAFÍA .....</b>	<b>503</b>
<b>ÍNDICE ANALÍTICO .....</b>	<b>533</b>

---

# Prefacio

---

No cabe la menor duda de que los conceptos abordados en el libro constituyen una sólida base de los avances, ya en marcha y de futuro, en Inteligencia Artificial donde el Aprendizaje Profundo es uno de sus máximos exponentes actuales.

Son muchos los ámbitos docentes, profesionales e industriales donde el término Aprendizaje Profundo suscita un inusitado interés, refiriéndose a él, la mayoría de las veces, como *Deep Learning* en su acepción inglesa.

Es evidente que los avances conceptuales y, sin duda, los tecnológicos han contribuido a dicho interés y desarrollo. Naturalmente, soportado por el buen desempeño que la puesta en práctica de los conceptos proporciona una vez instalados en múltiples y diversos dispositivos, que incorporan sistemas automáticos con diferentes propósitos. Robots y sistemas inteligentes en general, aplicaciones en telefonía móvil, reconocedores automáticos, entre otros, son buenos ejemplos de ello. Por otra parte, es necesario reconocer al respecto la contribución realizada por el desarrollo de tecnologías y modelos dentro del paradigma conocido como Internet de las Cosas (IoT, *Internet of Things*).

La base del Aprendizaje en general y del Profundo, en particular, es el procesamiento de los datos provenientes de diferentes sensores o de distintas fuentes de información. Ese procesamiento, en lo que respecta al Aprendizaje Profundo, se lleva a cabo mediante la aplicación de los conceptos abordados en el libro. Son innumerables los avances a nivel internacional al respecto, que han propiciado la aparición de ingentes publicaciones, y desarrollos en múltiples ámbitos y bajo todos los formatos posibles, digitales y convencionales. A este respecto, existe un problema evidente, que es el hecho de la dispersión de los conceptos bajo las diferentes apariencias mencionadas. Este libro trata de unir y aglutinar el mayor número posible de ellos bajo un mismo contenedor. Se facilita así, la posibilidad de iniciación o profundización en la materia, por parte de quienes tienen interés en ella,

tanto a nivel docente e investigador como profesional en el sector industrial. Es justamente aquí donde radica el valor añadido del libro, que proporciona, con suficiente claridad, los conceptos esenciales en sendos ámbitos con tal finalidad. Sin duda, esta es la razón fundamental por la que se publica la obra, siendo conscientes de la necesidad de un texto con las características reunidas por este. En este sentido, además, el libro se ha planteado de forma que resulte autosuficiente, de suerte que el lector, incluso sin ser experto en los conceptos expuestos, puede abordar sin dificultad los contenidos en el mismo. Sin duda, este hecho otorga un valor añadido de interés a la obra.

En definitiva, desarrolladores, ingenieros, investigadores o estudiantes universitarios encontrarán en el libro una referencia de base de suma utilidad para abordar los aspectos conceptuales y de implementación en el desarrollo de aplicaciones basadas en Aprendizaje Profundo, particularmente para quienes se inicien en la materia por su carácter didáctico y autocontenido.

## **Sobre los autores**

---

Gonzalo Pajares Martinsanz es profesor en la Facultad de Informática de la Universidad Complutense de Madrid, en el Departamento de Ingeniería del Software e Inteligencia Artificial y miembro del Instituto de Tecnología del Conocimiento de la misma Universidad. Ha desarrollado una extensa actividad profesional durante más de tres lustros en la industria aplicando tecnologías software y de Inteligencia Artificial. Además, está ampliamente involucrado en tareas de investigación en el ámbito de la Inteligencia Artificial mediante distintos proyectos nacionales e internacionales de investigación, durante más de dos décadas, con transferencia tecnológica a la industria. Es autor y editor de varios libros sobre visión por computador, inteligencia artificial, tecnologías sensoriales, así como autor de numerosas publicaciones en dichas áreas, incluyendo técnicas de reconocimiento de patrones y estructuras en varios ámbitos industriales, con participación de empresas situadas en la vanguardia tecnológica. Ha publicado numerosos artículos en revistas especializadas de prestigio internacional, a la vez que ha colaborado y colabora como editor invitado y asociado en varias revistas con alto índice de impacto donde las técnicas de Aprendizaje Profundo constituyen un aporte importante.

Pedro Javier Herrera Caro es profesor en la Escuela Técnica Superior de Ingeniería Informática de la Universidad Nacional de Educación a Distancia, en el Departamento de Ingeniería de Software y Sistemas Informáticos. Ha desarrollado su actividad investigadora y docente en la Universidad Complutense de Madrid, el Centro de Automática y Robótica (CSIC-UPM) y la Universidad Francisco de Vitoria. Ha realizado estancias de investigación en el Centro de Investigación Forestal del XII



Instituto Nacional de Investigación y Tecnología Agraria y Alimentaria, y en el AutoNOMOS LAB de la Universidad Libre de Berlín. Ha colaborado con diversos organismos y centros de investigación, universidades y empresas a través de proyectos internacionales, nacionales y autonómicos, y contratos con empresas. Fruto de ello son numerosas publicaciones en revistas, libros y conferencias internacionales. Su investigación se orienta hacia la Visión por Computador, el Reconocimiento de Patrones, la Inteligencia Artificial y la Robótica. Aplicando técnicas relacionadas con el procesamiento de imágenes, la clasificación de patrones, toma de decisiones, visión estereoscópica, teledetección y minería de datos.

Eva Besada Portas es profesora de la Sección Departamental de Arquitectura de Computadores y Automática de la Facultad de Ciencias Físicas de la Universidad Complutense de Madrid. Ha desarrollado una extensa actividad docente e investigadora durante los últimos veinte años, esta última centrada en el desarrollo y aplicación de técnicas de inteligencia artificial, optimización y fusión multisensorial a diferentes problemas reales. Es autora de numerosas publicaciones en estos ámbitos en revistas especializadas y congresos de prestigio internacional, y ha participado durante todos estos años en proyectos de investigación y contratos de transferencia tecnológica a la empresa relacionados con las mismas. Finalmente, indicar que ha realizado estancias de investigación tanto predoctorales como post-doctorales en centros extranjeros, como la Universidad de Londres, la Universidad Carnegie Mellon, la Universidad de Nuevo México o Google.

## Agradecimientos

---

Es digno reconocer la contribución activa o pasiva a todas las personas que han tenido que sufrir la sustracción de tiempo por culpa de este libro, familiares cercanos son los más afectados y se sitúan en primera línea. Un recuerdo especial a quienes siempre disfrutaron con el esfuerzo de sus allegados.

Gracias también a las industrias, por su confianza depositada en la aplicación con éxito de algunos conceptos en sus productos de mercado. Mención especial merece Lector Vision, como empresa tecnológica puntera en el procesamiento de imágenes. Gracias, también, a los respectivos Departamentos y al Instituto de Tecnología del Conocimiento de la Universidad Complutense de Madrid por la oportunidad brindada para la aplicación de los conceptos vanguardistas del libro en diversos ámbitos.

Finalmente, un especial agradecimiento a la editorial RC Libros por su soporte para la publicación y difusión de la obra. Gracias a José Luis, por su dedicación e incansable apoyo continuo a la difusión de los conceptos técnicos, en este caso de vanguardia tecnológica.

# 1 INTRODUCCIÓN

## 1.1 PRELIMINARES

---

El Aprendizaje Profundo constituye una rama muy poderosa del Aprendizaje Automático (AA), enmarcándose dentro de la Inteligencia Artificial (IA) y abriendo un campo de conocimiento muy prometedor en continuo auge, en parte gracias a los avances tecnológicos, que permiten el procesamiento de ingentes cantidades de datos con estructuras complejas, cuyo exponente máximo son las redes neuronales, y más específicamente las catalogadas como *profundas*. Es evidente que hoy en día la IA constituye un campo prometedor con múltiples y diversas aplicaciones prácticas, manteniendo un área de investigación muy activa, en la que han florecido abundantes sistemas inteligentes que facilitan y permiten automatizar el trabajo rutinario. Dentro de la IA destaca un área de gran relevancia, cual es el AA (Pajares y Cruz, 2010), que comprende una amplia gama de métodos y aplicaciones capaces de procesar ingentes cantidades de datos para extraer la máxima información posible subyacente en los mismos. En términos generales, el esquema global de aprendizaje consta de tres módulos principales, a saber:

a) *Generador*, que proporciona las entradas convenientemente estructuradas para su procesamiento, generalmente en forma de vectores como características, cuyas componentes son los atributos de los datos; por ejemplo, si los datos son imágenes, las características pueden ser los píxeles y los atributos los valores de color de estos.

b) *Entrenamiento*, de forma que a cada vector se le asigna una categoría de salida, constituyendo el objetivo para ajustar una serie de parámetros, que son el resultado del aprendizaje, finalizando con el modelo correspondiente convenientemente estructurado o aprendido. La mencionada asignación puede ser de naturaleza supervisada o no supervisada, en función de que haya un instructor o no, respectivamente, dirigiendo el proceso. Los vectores de entrada se denominan patrones (*patterns*), casos (*cases*), entradas (*inputs*), instancias (*instances*) u observaciones (*observations*). Las categorías de salida se denominan etiquetas (*labels*), objetivos (*targets*), salidas (*outputs*) y a veces también observaciones (*observations*).

c) *Decisión*, que asigna una categoría dada a una nueva muestra de entrada aplicando los parámetros aprendidos del modelo. La *clasificación*, la *regresión* o la *optimización*, entre otras, son tareas que pertenecen a esta fase.

Dentro del AA se encuadra el conocido como *Aprendizaje Profundo* (AP). En la literatura especializada a nivel internacional es muy común referirse al AP por sus términos en inglés, esto es, *Deep Learning*, que es el núcleo central del presente libro. Por otra parte, y al hilo de esta cuestión, conviene reseñar que son muchos y diversos los términos en inglés utilizados para definir y describir los conceptos involucrados bajo el paradigma del AP. Muchos de los cuales no poseen una clara traducción conceptual al español, razón por la cual los considerados bajo esta situación se mantienen a lo largo del libro con el fin de que el lector pueda fácilmente identificarlos en la literatura especializada escrita en inglés. Solo se han traducido aquellos conceptos que no admiten discusión, manteniendo en todo caso su expresión original en inglés.

Bien es cierto que desde los años 50 del siglo pasado, la IA a veces se ha sobrevalorado y se ha considerado como muy prometedora en diversas ocasiones, eso pese a que no se han llegado a alcanzar las perspectivas iniciales. Por otra parte, no es menos cierto que en los últimos años se están viendo avances importantes gracias al AP. Ello a pesar de que todavía es relativamente frágil de cara a su generalización y adaptación a entornos o escenarios cambiantes, principalmente por falta de datos suficientes que capturen los cambios de dicho entorno, pudiendo aparecer ciertos sesgos por falta de la información necesaria extraíble de los datos.

Algunos autores como Marcus y Davis (2019), expresan, no sin cierta razón, algunos aspectos relativos a las ventajas e inconvenientes de los procesos de AP, achacándoles que muchas arquitecturas basadas en redes neuronales hacen cosas increíbles, sin ser conscientes por parte de quien las aplica del conocimiento real sobre lo que están haciendo, y por ello, no son sistemas totalmente inteligentes.

Aunque en parte, esto último puede verse desde esta perspectiva, no es menos cierto que los desarrollos basados en AP son capaces de conseguir resultados importantes, siendo esta la perspectiva desde la que se abordan y plantean los temas del presente libro.

Como sostienen Goodfellow y col. (2016), en los comienzos de la IA, se abordaron rápidamente problemas intelectualmente difíciles para los seres humanos, pero relativamente sencillos para las computadoras, todo ello mediante una lista de reglas matemáticas formales. A partir de ahí, el verdadero desafío para la IA se tornó en resolver tareas fáciles de realizar para las personas, pero difíciles de describir formalmente. Aquí se incluyen tareas tales como el reconocimiento de objetos en imágenes, palabras o acciones en los movimientos. No cabe duda de que en este aspecto el AP ha conseguido ya logros muy relevantes. Es en este rasgo, y más concretamente en la exposición de una serie de técnicas orientadas a tal fin, donde se centra el presente libro.

En definitiva, se trata de exponer una serie de técnicas para resolver problemas, por decirlo de alguna manera, intuitivos para el ser humano, con el uso de las computadoras, que aprenden mediante los métodos y algoritmos diseñados a partir de los datos suministrados y sin necesidad de que los humanos especifiquen formalmente todo el conocimiento requerido por la computadora. En cualquier caso, y siguiendo también la teoría expuesta en Goodfellow y col. (2016), la jerarquía de conceptos permite que la computadora aprenda conceptos complicados al construirlos a partir de otros más simples, todos ellos estructurados en múltiples capas, razón por la cual a este enfoque se le denomina con el término ya indicado de AP. En cualquier caso, como una técnica específica del AA, estos procedimientos se encaminan a extraer patrones determinados a partir de los datos.

Existe una diferencia fundamental en lo que respecta a la extracción de las mencionadas características entre las técnicas clásicas, por llamarlas de alguna manera, de AA y las específicas del AP. Por ejemplo, considérese un ejemplo sencillo biclase, en el que se trata de separar el cielo y la hierba en una imagen de color de un paisaje de campo. Los datos disponibles en este caso son valores de color de los píxeles, de forma que en el caso del cielo predominan las tonalidades azules, mientras que en la hierba son las verdes. Un método simple tal como *naive Bayes* puede separar los patrones en dos clases diferentes teniendo en cuenta que los mismos están definidos, por lo que se conoce como características. La extracción de características en este caso es esencial. Por otro lado, y siguiendo en el ámbito de las imágenes, estas se caracterizan por poseer información espacial, y en vídeos, también temporal. Las *redes neuronales profundas* pueden captar perfectamente ambos tipos de información. En el primer caso, los *filtros de convolución* son

responsables de la captura espacial, a diferencia de lo que ocurre con otros modelos de red, tal como las de *retropropagación*, en las que las características de las imágenes se transforman en vectores que se suministran a la entrada, perdiendo las relaciones espaciales. Por ejemplo, una imagen de dimensión 3x3 se transforma en un vector con 9 componentes. En el caso de las características temporales las redes recurrentes tienen la habilidad de realizar tal captura.

No obstante, para muchas tareas no resulta fácil extraer características. Por ejemplo, supóngase que a partir de una imagen se quieren identificar peatones cuando un vehículo autónomo navega en un entorno urbano. Una persona puede identificarse por poseer cabeza, tronco y extremidades. Se podría pensar en detectar la presencia de extremidades o del cuerpo y la cabeza o todas, lo cual no resulta trivial debido a que no es fácil establecer las características de dichas partes. Los brazos y las piernas son alargados, el tronco tiene una forma más rectangular, pero en todos los casos nunca están exentos de elementos perturbadores, tales como el uso de distintos tipos de ropa, sombras, oclusiones totales o parciales, entre otros.

Una solución a este problema consiste en utilizar el AA para descubrir no solo la proyección de la representación a la salida, sino también la representación misma. Este enfoque se conoce como *aprendizaje de representación* según Goodfellow y col. (2016). Las representaciones aprendidas a menudo resultan en un rendimiento mucho mejor que el que se puede obtener con representaciones diseñadas a mano. También permiten que los sistemas de IA se adapten rápidamente a nuevas tareas, con una mínima intervención humana. Un algoritmo de aprendizaje de representación puede descubrir un buen conjunto de características para una tarea simple en minutos o para una tarea compleja en horas y meses. El diseño manual de características para una tarea compleja requiere una gran cantidad de tiempo y esfuerzo humanos, pudiendo llevar incluso décadas. Un ejemplo por excelencia de un algoritmo de aprendizaje de representación es el *autoencoder* (autocodificador), que convierte los datos de entrada en una representación diferente para luego poder devolverla a la representación original mediante el correspondiente *decoder* (decodificador).

Cuando se diseñan algoritmos para aprender las características, el objetivo consiste en separar los factores de variación que expliquen los datos observados. El concepto *factor* se refiere a abstracciones que ayudan a distinguir entre la alta variabilidad de los datos observados; así, en el reconocimiento de los peatones los factores de variación hacen referencia a la posición de las extremidades con respecto al tronco, la posición con respecto a la cámara, la ropa con la que van vestidos, las oclusiones de las extremidades, las posibles sombras proyectadas sobre sus cuerpos o la intensidad de la luz con la que se ha obtenido la imagen, entre otros. La mayoría

de las aplicaciones exigen separar los factores de variación descartando aquellos que no interesan. A la vista de lo cual, resulta francamente difícil obtener una representación tal que permita resolver el problema. Es precisamente aquí donde entra en acción el aprendizaje profundo, ya que permite introducir representaciones que se expresan en términos de otras representaciones a distintos niveles, que van estructurando convenientemente la información. Por ejemplo, la figura 1-1 muestra un sistema basado en AP, concretamente una Red Neuronal Convolutiva (RNC) o en terminología inglesa *Convolutional Neural Network* (CNN), donde se representa el concepto de una imagen de una taza combinando conceptos más simples, tales como bordes, contornos o partes de los objetos hasta llegar a su clasificación, en este caso como taza. La idea de representación en múltiples capas es lo que determina una de las perspectivas del AP. De esta forma, se puede decir con carácter general, que las primeras capas de las redes profundas extraen características de bajo nivel, de modo que estas se van tornando en más complejas con características de mayor nivel hasta llegar a las capas superiores, en las que las características extraídas de la imagen son del más alto nivel. Esta información concatenada permite identificar un objeto (taza), a pesar de que pueda presentar diferentes características tales como, por ejemplo, color, forma o tamaño, lo que permite claramente diferenciar el AP del AA.

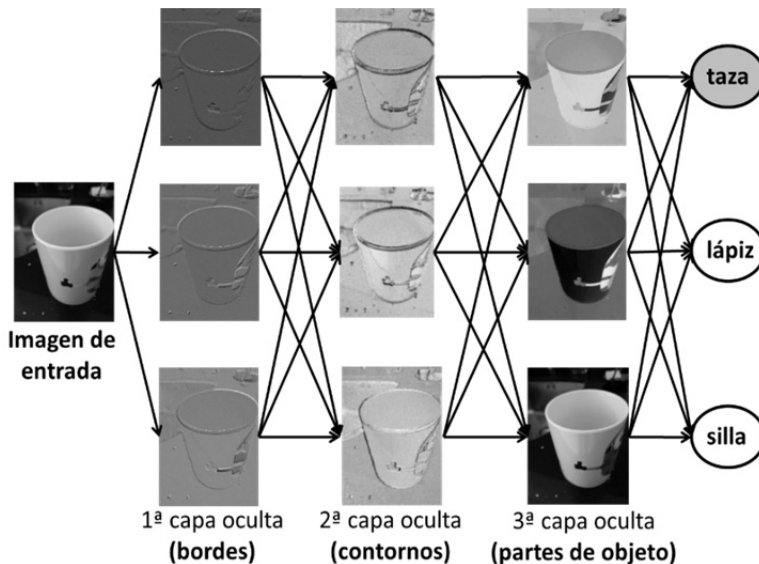


Fig. 1-1 Modelo de aprendizaje profundo

Otra idea para determinar el concepto de *profundidad* es la también establecida por Goodfellow y col. (2016), en el sentido de que la profundidad se determina como el estado del computador para aprender un programa computacional multipaso, de forma que cada capa de la representación puede verse como el estado de la

memoria del computador después de ejecutar otro conjunto de instrucciones en paralelo. Las redes con mayor profundidad pueden ejecutar más instrucciones en secuencia. Las instrucciones secuenciales ofrecen un gran poder porque las instrucciones posteriores pueden referirse a los resultados de instrucciones anteriores. Según esta visión del AP, no toda la información en las activaciones de una capa codifica necesariamente factores de variación que explican la entrada. La representación también almacena información de estado que ayuda a ejecutar un programa que puede dar sentido a la entrada. Esta información de estado podría ser análoga a un contador o puntero en un programa de computación tradicional. No tiene nada que ver con el contenido de la entrada específicamente, pero ayuda al modelo a organizar su procesamiento. Existen dos formas principales de medir la profundidad de un modelo. La primera se basa en el número de instrucciones secuenciales que deben ejecutarse para evaluar la arquitectura. Se puede pensar en esto como la longitud de la ruta más larga a través del diagrama de flujo que describe cómo calcular cada una de las salidas del modelo dadas sus entradas. Otro enfoque, utilizado por modelos probabilísticos profundos, considera que la profundidad de un modelo no es la profundidad del gráfico computacional sino la profundidad del gráfico que describe cómo se relacionan los conceptos entre sí.

En este caso, la profundidad del diagrama de flujo de los cálculos necesarios para computar la representación de cada concepto puede ser mucho más profunda que la gráfica de los conceptos en sí mismos. Esto se debe a que la comprensión del sistema de los conceptos más simples puede refinarse dando información sobre los conceptos más complejos. Por ejemplo, siguiendo también a Goodfellow y col. (2016), un sistema inteligente que observa una imagen de una cara con un ojo en la sombra puede ver inicialmente únicamente un ojo. Después de detectar la presencia de una cara, el sistema puede inferir que probablemente también esté presente un segundo ojo. En este caso, la gráfica de conceptos incluye sólo dos capas, una capa para *ojos* y una capa para *caras*, pero la gráfica de cálculos incluye dos capas si se refina la estimación de cada concepto dadas las otras  $n$  veces. Debido a que no siempre está claro cuál de estos dos modelos (la profundidad del gráfico computacional o la profundidad del gráfico de modelado probabilístico) es más relevante, y debido a que distintas personas eligen diferentes conjuntos de elementos más pequeños a partir de los cuales construir sus gráficos, no existe un único valor correcto para la profundidad de una arquitectura, así como tampoco hay un valor correcto único para la duración de un programa computacional. Tampoco existe un consenso acerca de la profundidad que un modelo requiere para calificarlo como “profundo”. Sin embargo, el aprendizaje profundo puede considerarse, con seguridad, como el estudio de modelos que implican una mayor cantidad de composición de funciones aprendidas o conceptos aprendidos que el aprendizaje automático tradicional.

En definitiva, el aprendizaje profundo es un tipo particular de aprendizaje automático que consigue un gran poder y flexibilidad al representar al mundo como una jerarquía anidada de conceptos, donde cada concepto se define en relación con conceptos más simples y representaciones más abstractas calculadas en términos de conceptos menos abstractos.

## 1.2 BREVE REVISIÓN HISTÓRICA

---

Tal y como sostienen Goodfellow y col. (2016), el concepto de aprendizaje profundo (AP) data de los años 40 del siglo pasado a pesar de su actual relevancia y auge, lo que ciertamente parece sorprendente. Ello se debe al hecho de que durante cierto tiempo las aplicaciones basadas en este paradigma no consiguieron el éxito que se les suponía inicialmente.

Una primera etapa corresponde a la época en la que el AP se enmarca en el contexto cibernético, aproximadamente entre los años 1940 y 1960 y coincidiendo con la época del desarrollo de las teorías biológicas del aprendizaje como las desarrolladas por McCulloch y Pitts (1943) o Hebb (1949) a las que siguieron las implementaciones de los primeros modelos, entre los que destaca el perceptrón (Rosenblatt, 1958), que es la unidad básica en las redes neuronales o ADALINE (Widrow y Hoff, 1960). Por tanto, en esta fase prevalecen los modelos lineales, de forma que a partir de un conjunto de  $n$  valores de entrada, la salida es una combinación lineal en la que intervienen los pesos de un modelo que se ajustan durante la fase de aprendizaje. No obstante, la incapacidad para resolver el problema XOR con modelos lineales llevó, en parte, a la decadencia en la popularidad de las redes neuronales.

Las redes neuronales artificiales fueron uno de los nombres que recibió el AP, evolucionando posteriormente a un principio más general de aprendizaje en el sentido de configurarse como múltiples niveles de composición. No obstante, a diferencia de aquellos modelos iniciales basados en la neurociencia, ahora en el AP la neurociencia es la inspiración, pero no el modelo, entre otras razones por el desconocimiento del cerebro humano para considerarse una guía, ya que incluso se desconocen operaciones del cerebro elementales (Olshausen y Field, 2005).

No obstante, la neurociencia proporciona una perspectiva avanzada en el ámbito del AP proporcionando la idea de que quizás un único algoritmo en este campo puede ser válido para resolver distintas tareas, tal y como sostienen Goodfellow y col. (2016). Esto es así a raíz del descubrimiento por parte de Von Melchner y col. (2000), de que los hurones pueden llegar a ver cosas con la región auditiva de su cerebro. Esto ha permitido que diversas disciplinas converjan hacia desarrollos y