

# Los retos del aprendizaje máquina en la era del Big Data

## Machine Learning Challenges in Big Data Era

Miguel Veganzones-Bodón  
Sherpa.ai (España)

DOI: <http://dx.doi.org/10.6036/9243>

### INTRODUCCIÓN

Durante la última década hemos asistido a la consolidación de una serie de tecnologías disruptivas, como el Internet de las Cosas (*Internet of Things*) y la computación escalable (*Big Data*), así como a la popularización del análisis avanzado de datos (*Data Science*) y otras técnicas asociadas a la Inteligencia Artificial, por ejemplo, el Aprendizaje Máquina (*Machine Learning*), la Visión Artificial o los Sistemas Conversacionales. Muchas grandes compañías tecnológicas han fundamentado su éxito en estos avances, de los cuales han sido, en gran medida, máximas precursoras y adalides. Reflejado en este éxito y dada la ubicua naturaleza de los datos, el mercado global busca replicarlo en todo tipo de industrias, en parte por la oportunidad de crear nuevo valor y nego-

cios a partir de los datos, y en parte, por temor a quedarse atrás.

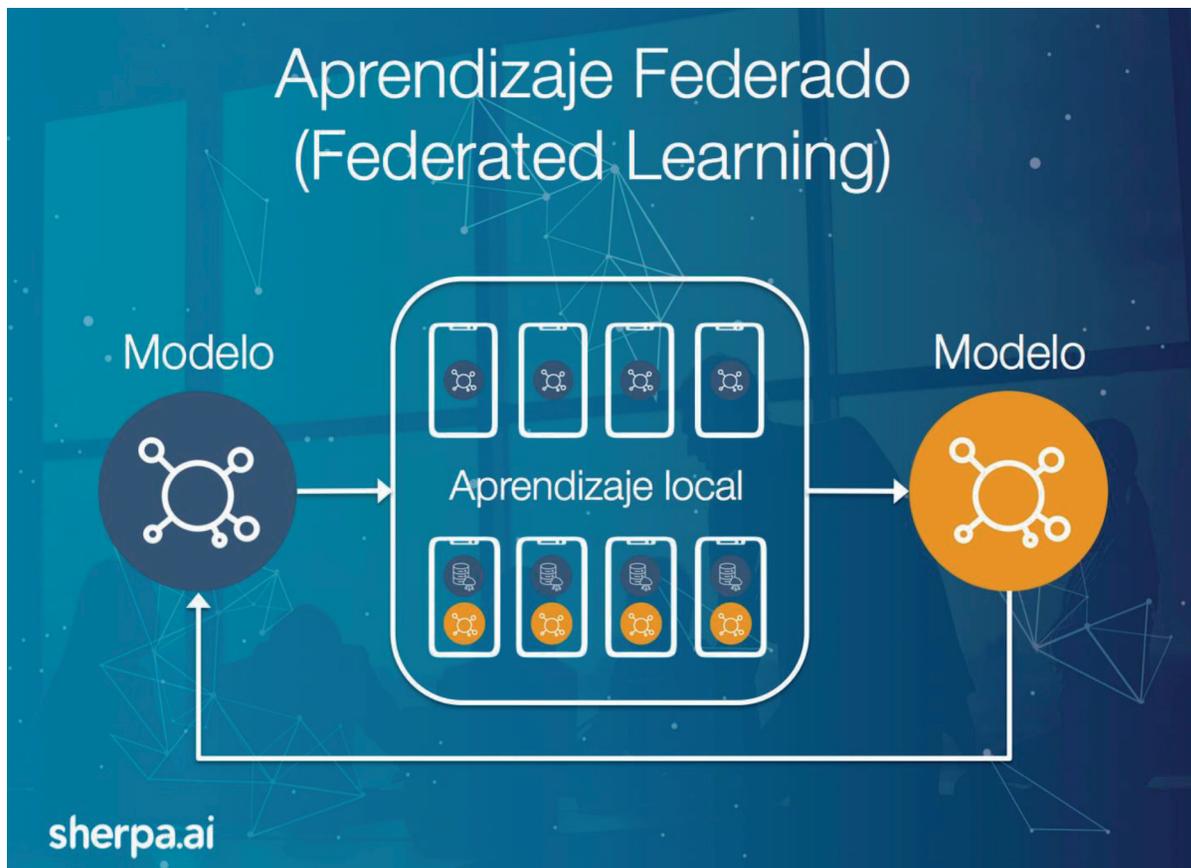
Esta situación es sintomática de un cambio de paradigma cuyo eco no se limita a modificar los tejidos científico-tecnológico e industrial. La sociedad en su conjunto ha tomado un rol activo en este cambio, tanto en las comunidades profesionales como en nuestro día a día. Continuamente, generamos ingentes cantidades de datos que facilitan el desarrollo de tecnologías que nos apoyan en la decisión, o que directamente deciden por nosotros en aspectos relevantes de nuestras vidas. Tecnologías que adaptan la oferta de ocio y consumo a nuestros gustos y necesidades, pero también tecnologías que deciden sobre nuestra salud, finanzas, educación y otros aspectos fundamentales que nos definen como ciudadanos con derechos y obligaciones.

El Prof. Thomas Kuhn expuso en su obra de 1962 "La estructura de las revoluciones científicas" [1] que las ciencias naturales evolucionan en base a alternancias de períodos de "ciencia normal" y cambios

abruptos. Según Kuhn, la "ciencia normal" se define mediante una serie de creencias tácitas sobre los valores, la dirección y los objetivos de la investigación científica. Además, indica que la comunidad científica toma como referentes ejemplos paradigmáticos que describen cómo se debe hacer la ciencia. Un cambio de paradigma es un cambio abrupto en lo que definimos como "ciencia normal" y, por ende, en los paradigmas de referencia.

### EVOLUCIÓN FUTURA

La comunidad de la Inteligencia Artificial, y particularmente del Aprendizaje Máquina, ha estado inmersa en un cambio de paradigma que ha trascendido a la propia comunidad científica [2]. El paradigma anterior se fundamentaba en el desarrollo empírico y teórico de métodos de aprendizaje a partir de datos y recursos limitados, apoyado por modelos con fuertes asunciones implícitas y/o explícitas. El valor fundamental de la Inteligencia Artificial residía en entender cómo estos





modelos podían generalizar lo que habían aprendido de unos pocos datos curados en ciertos casos paradigmáticos, como el reconocimiento de dígitos y caras, el test de Turing o la práctica de juegos complejos como el ajedrez o el go.

El nuevo paradigma en el campo de la Inteligencia Artificial está fuertemente marcado por la disposición de ingentes cantidades de datos heterogéneos, a menudo recogidos en tiempo real, y los avances tecnológicos que permiten escalar linealmente en recursos para procesar y analizar tal cantidad de datos. La principal característica de este cambio de paradigma es que muchos de los problemas clásicos de la Inteligencia Artificial han sido resueltos durante la última década generando grandes cantidades de datos y entrenando sistemas para directamente emular el comportamiento deseado, en lugar de tratar de emular la forma en que dicho comportamiento se aprende [3]. Esta es una forma radicalmente distinta de enfocar el Aprendizaje Máquina, que pone el foco en la capacidad predictiva de los datos en lugar de en la comprensión de lo aprendido por los modelos.

Este desplazamiento de la carga del aprendizaje hacia los datos ha favorecido el desarrollo de modelos que ya eran conocidos en la comunidad científica pero que permanecían en un segundo plano. Así, el cambio de paradigma ha conllevado el desplazamiento de recursos hacia la investigación de modelos de aprendizaje profundo (Deep Learning) [4], que se han mostrado superiores a la hora de resolver problemas específicos del campo, como en los avances recientes para Sistemas Conversacionales y Asistentes Inteligentes,

únicamente a partir de grandes volúmenes de datos, sin requerir curar los datos ni incluir conocimiento del dominio en los modelos, procesos que suelen ser costosos y ofrecen muchas dificultades operacionales.

La tremenda importancia de los datos da pie a una mayor preocupación por los aspectos éticos y legales del uso de la información sensible recogida en ellos. La Privacidad de los Datos se ha convertido en uno de los principales temas de discusión en la comunidad, tratado de formas dispares en los marcos jurídicos de los distintos organismos legisladores. Sin duda alguna es una dificultad que puede poner en serio riesgo la capacidad de trasladar a la ciudadanía y la industria los avances en el campo. Por ello, una línea de trabajo que está cobrando cada vez mayor importancia es el Aprendizaje Federado (Federated Learning) [5], marco en el que se estudian mecanismos para entrenar modelos sin necesidad de acceder directamente a los datos, sino en base a agregar mejoras aprendidas de forma privada, por ejemplo, en los dispositivos móviles de los usuarios. De esta manera, los individuos darían acceso a las mejoras incrementales de los modelos pero no a sus datos, manteniendo la privacidad de éstos.

La Interpretabilidad de los Modelos es otro objeto de discusión y estudio abundante. Desplazar el foco hacia la capacidad predictiva de los modelos intensivos de datos, hace que éstos se consideren prácticamente como cajas negras, y sea muy difícil interpretar las decisiones tomadas por una IA. Cuando el objeto de la IA es ofrecer recomendaciones personales sobre el ocio y el consumo, éste puede pa-

recer un mal menor, pero es de vital importancia cuando estas decisiones afectan derechos fundamentales de los ciudadanos, que pueden ser socavados por sesgos humanos recogidos en los datos y aprendidos por los modelos, por ejemplo, sesgos raciales o de género. Nuevas técnicas para tratar de identificar y reducir estos sesgos, así como para interpretar las decisiones de modelos complejos, tales como los modelos de Aprendizaje Profundo, resultarán fundamentales para de nuevo, hacer que estas Inteligencias Artificiales sean viables.

Finalmente, creo oportuno resaltar una vía de investigación que puede ser el nuevo gran frente abierto del campo de la IA: humanizar las Inteligencias Artificiales del futuro, dotándolas de cierta capacidad empática para responder a los sentimientos y estados emocionales de las personas. Las Inteligencias Artificiales Empáticas [6] deberán ser capaces de adaptarse a situaciones complejas, donde el contexto emocional es fundamental para la toma de decisiones.

## REFERENCIAS

- [1] Kuhn, Thomas S. "La estructura de las revoluciones científicas", Fondo de Cultura Económica de España, ISBN 978-84-375-0579-4.
- [2] Cristianini, Nello. "On the current paradigm in artificial intelligence", *AI Communications*, vol. 27, no. 1, pp. 37-43, 2014.
- [3] Halevy, Alon; Norvig, Peter; Pereira, Fernando. "The Unreasonable Effectiveness of Data", *IEEE Intelligent Systems*, vol. 24, pp. 8-12, 2009.
- [4] Goodfellow, Ian; Bengio, Yoshua; Courville, Aaron. "Deep Learning", MIT Press, 2016.
- [5] <https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>
- [6] Mäntylä, Mika V.; Graziotin, Daniel; Kuutila, Miikka. "The evolution of sentiment analysis - A review of research topics, venues, and top cited papers", *Computer Science Review*, vol. 27, pp. 16-32, 2018.