

ERIK J. LARSON

**EL MITO
DE LA
INTELIGENCIA
ARTIFICIAL**

*Por qué las máquinas
no pueden pensar como
nosotros lo hacemos*

EL MITO DE LA INTELIGENCIA ARTIFICIAL

EL MITO DE LA INTELIGENCIA ARTIFICIAL

Por qué las máquinas no pueden
pensar como nosotros lo hacemos

ERIK J. LARSON

Traducción de Milo J. Krmpotić

Shackleton
— b o o k s —

El mito de la Inteligencia Artificial. Por qué las máquinas no pueden pensar como nosotros lo hacemos

Título original: *The Myth of Artificial Intelligence. Why computers can't think the way we do*

© 2021 Erik J. Larson

© de esta edición, Shackleton Books, S. L., 2022

La presente edición se publica en acuerdo con Harvard University Press a través de International Editors' Co

© Traducción: Milo J. Krmptic

Shackleton
— b o o k s —



@Shackletonbooks

www.shackletonbooks.com

Realización editorial: La Letra, S. L.

Diseño de cubierta: Pau Taverna

Conversión a ebook: Iglú ebooks

ISBN: 978-84-1361-202-7

Reservados todos los derechos. Queda rigurosamente prohibida la reproducción total o parcial de esta obra por cualquier medio o procedimiento y su distribución mediante alquiler o préstamo públicos.

ÍNDICE

Introducción

Primera parte. El mundo simplificado

Capítulo 1. El error de la inteligencia

Capítulo 2. Turing en Bletchley.

Capítulo 3. El error de la superinteligencia

Capítulo 4. La singularidad, ayer y hoy.

Capítulo 5. La comprensión del lenguaje natural

Capítulo 6. De la IA como tecnología kitsch

Capítulo 7. Simplificaciones y misterios

Segunda parte. El problema de la inferencia

Capítulo 8. No calcules, analiza

Capítulo 9. El puzle de Peirce (y el rompecabezas de Peirce)

Capítulo 10. Problemas de deducción e inducción

Capítulo 11. El aprendizaje automático y el big data

Capítulo 12. La inferencia abductiva

Capítulo 13. Inferencia y lenguaje 1

Capítulo 14. Inferencia y lenguaje 2

Tercera parte. El futuro del mito

Capítulo 15. Mitos y héroes

Capítulo 16. La mitología de la IA invade la neurociencia

Capítulo 17. Las teorías de la inteligencia humana basadas en el neocórtex

Capítulo 18. ¿El fin de la ciencia?

Agradecimientos

Notas

Para Brooke y Ben

Introducción

En las páginas de este libro vas a leer acerca del mito de la inteligencia artificial. Lo de «mito» no se refiere a la imposibilidad de una IA verdadera. A ese respecto, el futuro de la IA es un misterio para la ciencia. El mito de la inteligencia artificial consiste en afirmar que su llegada es inevitable, mera cuestión de tiempo —que nos hemos adentrado ya en el sendero que conducirá a una IA de nivel humano, y más tarde a una superinteligencia—. No es así. Ese sendero existe solo en nuestra imaginación. Sin embargo, el carácter inevitable de la IA se encuentra tan arraigado en el debate popular —promovido por los expertos de los medios de comunicación, por referentes intelectuales como Elon Musk e incluso por numerosos científicos de IA (aunque desde luego no por todos ellos)— que, a menudo, cualquier pega que se le ponga se considera una forma de ludismo, o por lo menos una visión corta de miras sobre el futuro de la tecnología y un fracaso peligroso a la hora de prepararse para un mundo de máquinas inteligentes.

Tal y como os voy a mostrar, la ciencia de la IA ha revelado un misterio de grandes dimensiones en el núcleo de la inteligencia, y en la actualidad nadie tiene la menor idea de cómo resolverlo. Los partidarios de la IA cuentan con inmensos incentivos para minimizar sus limitaciones. Al fin y al cabo, la IA es un negocio enorme y tiene una presencia cada vez más predominante en la cultura. No obstante, nos guste o no, la posibilidad de un futuro de sistemas de IA se encuentra limitada por lo que sabemos en la actualidad sobre la naturaleza de la inteligencia. Y aquí deberíamos afirmarlo con franqueza: todas las pruebas sugieren que las inteligencias humana y artificial son radicalmente diferentes. El mito de la IA insiste en

que esas diferencias son solo temporales, y en que la aparición de sistemas más potentes acabará por erradicarlas. Futurólogos como Ray Kurzweil y el filósofo Nick Bostrom, prominentes proveedores del mito, hablan no solo como si la IA de nivel humano resultara inevitable, sino como si, al poco de su llegada, las máquinas superinteligentes fueran a dejarnos muy atrás.

Este libro explica dos aspectos importantes del mito de la IA, uno de tipo científico y otro cultural. La parte científica del mito asume que solo tenemos que seguir «desnudando la cebolla» del desafío de la inteligencia general, avanzando en hitos restrictivos de la inteligencia como la participación en juegos o el reconocimiento de imágenes. Se trata de un error grave: el éxito en las aplicaciones débiles no nos acerca ni un solo paso a la inteligencia general. Las inferencias que requieren los sistemas de cara a alcanzar una inteligencia general —leer el periódico, o mantener una conversación elemental, o ejercer de ayudante, como el robot Rosie de *Los Supersónicos*— no se pueden programar, aprender ni diseñar a partir de nuestro conocimiento actual de la IA. Al aplicar con éxito versiones de inteligencia más simples y débiles, que se benefician del uso de ordenadores más rápidos y de montones de datos, no estamos obteniendo un avance progresivo, sino que nos limitamos a recoger sus frutos maduros. El salto hacia un «sentido común» general es completamente diferente, y no se conoce camino alguno que lleve de lo uno a lo otro. No existe ningún algoritmo para la inteligencia general. Y tenemos buenos motivos para mostrarnos escépticos ante la idea de que dicho algoritmo vaya a surgir de nuevas tentativas con los sistemas de aprendizaje profundo o de cualquier otra aproximación popular en la actualidad. Resulta mucho más probable que vaya a requerir de un avance científico de primer orden, y ahora mismo nadie tiene la más remota idea del aspecto que tendría ese avance, y mucho menos de los detalles que conducirán a él.

La mitología sobre la IA es negativa, pues, porque oculta un misterio científico bajo la cháchara interminable del progreso continuado. El mito sostiene la creencia en un éxito inevitable, pero el respeto genuino por la ciencia debería hacer que volviéramos a la casilla de salida. Eso nos conduce al segundo tema de estas páginas: las consecuencias culturales del mito. Perseguir un mito no es la mejor manera de obtener «inversiones expertas», y ni siquiera una posición neutral. Es malo para la ciencia y es malo para nosotros. ¿Por qué? Un motivo es que resulta poco probable que alcancemos innovaciones si decidimos ignorar un misterio tan básico en vez

de afrontarlo. La versión saludable de la cultura de la innovación pone el énfasis en la exploración de lo que se desconoce, no en dar bombo a la ampliación de unos métodos ya existentes —sobre todo cuando esos métodos se han revelado inadecuados para llevarnos mucho más allá—. La mitología acerca del éxito inevitable de la IA tiende a extinguir la cultura misma de la invención, tan necesaria para obtener un avance real —con la IA de nivel humano o sin ella—. El mito también fomenta la resignación ante el progresivo avance hacia una tierra de máquinas, donde la invención genuina se deja de lado en favor de charlas futuristas que defienden los métodos actuales, a menudo desde intereses particulares.

¿Quién debería leer este libro? Sin duda, cualquier persona que se emocione con la idea de la IA pero que se esté preguntando por qué siempre aparece a diez o veinte años vista. Hay un motivo científico para ello, que explicaré. También deberías leer este libro si piensas que el progreso de la IA hacia la superinteligencia es inevitable y te preocupa lo que habrá que hacer cuando llegue. Aunque no puedo demostrar que una *autoridad suprema* de la IA no vaya a aparecer algún día, sí puedo ofrecerte razones para que descartes con rigurosidad la perspectiva de ese escenario. Más en general, debes leer este libro si sientes curiosidad pero a la vez te encuentras confundido por el bombo generalizado que rodea a la IA en nuestra sociedad. Te explicaré los orígenes del mito de la IA, lo que sabemos y lo que ignoramos acerca de la perspectiva de alcanzar una IA de nivel humano, y el motivo por el que deberíamos apreciar mejor la única inteligencia verdadera que conocemos: la nuestra.

EN ESTE LIBRO

En la primera parte, «El mundo simplificado», explico que la cultura de la IA nos ha llevado a simplificar nuestras ideas sobre la gente a la vez que expandía nuestro conocimiento acerca de la tecnología. Esto comenzó con el fundador de la IA, Alan Turing, e incluye una serie de simplificaciones comprensibles pero desafortunadas que yo denomino «errores de inteligencia». Esos errores iniciales fueron magnificados hasta acabar conformando una ideología por parte de un amigo de Turing, el estadístico I. J. Good, quien introdujo la idea de «ultrainteligencia» como resultado

predecible tras la consecución de una IA de nivel humano. Entre Turing y Good vemos cobrar forma al mito moderno de la IA. Su desarrollo nos ha conducido a una época de lo que yo llamo «tecnología *kitsch*», imitaciones baratas de ideas más profundas que anulan el compromiso inteligente y debilitan nuestra cultura. Lo *kitsch* nos indica lo que hemos de pensar y lo que hemos de sentir. Los proveedores del *kitsch* sacan rédito de él, mientras que los consumidores de ese *kitsch* experimentan una pérdida; acaban —acabamos— metidos en un mundo de frivolidad.

En la segunda parte, «El problema de la inferencia», argumento que no tenemos la menor idea sobre cómo programar o diseñar el único tipo de inferencia —de pensamiento, en otras palabras— que funcionará con una IA de nivel humano (o cualquier otra cosa que se le acerque). El problema de la inferencia apunta al corazón del debate sobre la IA porque trata directamente con la inteligencia, la de la gente o la de las máquinas. Nuestro conocimiento acerca de los distintos tipos de inferencia se remonta a Aristóteles y a otros griegos de la Antigüedad, y se ha desarrollado en los ámbitos de la lógica y de las matemáticas. La inferencia ya se describe usando sistemas formales y simbólicos como los programas informáticos, así que explorándola se puede obtener una visión muy clara del proyecto con el que diseñar la inteligencia. Hay tres tipos de inferencia. La IA clásica exploró uno (las deducciones), la IA moderna explora otro (las inducciones). Y el tercer tipo (las abducciones) conduce a la inteligencia general y, sorpresa: nadie está trabajando en él —nadie en absoluto—.¹ Por último, puesto que todos los tipos de inferencia son distintos —con ello quiero decir que ninguno de esos tipos puede rebajarse hasta convertirse en otro—, sabemos que un fracaso a la hora de construir sistemas de IA que usen el tipo de inferencia en el que se afianza la inteligencia general conducirá al fracaso de los avances hacia la inteligencia artificial general, o IAG.

En la tercera parte, «El futuro del mito», argumento que, cuando se lo toma uno en serio, el mito tiene consecuencias muy negativas, ya que subvierte la ciencia. En especial, erosiona la cultura de la invención y la inteligencia humanas, que resultan necesarias en aquellos descubrimientos imprescindibles para comprender nuestro propio futuro. La ciencia de datos (la aplicación de la IA a los macrodatos) es, en el mejor de los casos, una prótesis del ingenio humano; en caso de usarla de manera correcta, nos ayudará a lidiar con el «diluvio de datos» contemporáneo. Cuando se la usa

para reemplazar la inteligencia individual, tiende a estropear la inversión sin ofrecer ningún resultado. Explico, en especial, que el mito ha afectado negativamente la investigación en neurociencia, entre otros avances científicos recientes. Estamos pagando un precio demasiado elevado por este mito. Como no poseemos ninguna buena razón científica para creer que el mito pueda hacerse realidad, puesto que contamos con todos los motivos para rechazarlo a fin de alcanzar la prosperidad en el futuro, tenemos que repensar de manera radical la conversación sobre la IA.

Primera parte
EL MUNDO SIMPLIFICADO

Capítulo 1

El error de la inteligencia

La historia de la inteligencia artificial comienza con las ideas de una persona que contó con una enorme inteligencia humana: el pionero de la informática Alan Turing.

En 1950, Turing publicó un artículo provocador, «Maquinaria computacional e inteligencia», sobre la posibilidad de crear máquinas inteligentes.¹ Fue un texto audaz, que llegó en un momento en el que los ordenadores eran novedosos pero insignificantes, según los parámetros de hoy en día. Aquellas piezas pesadas y lentas de *hardware* servían para acelerar cálculos científicos como el del análisis criptográfico. Tras una larga preparación, se les podían proporcionar fórmulas de física y unas condiciones iniciales, y obtener de forma automática el radio de una explosión nuclear. IBM no tardó en entender su potencial de cara a reemplazar a los seres humanos en sus operaciones comerciales, como la actualización de hojas de cálculo. Pero ver los ordenadores como criaturas «pensantes» requería de cierta imaginación.

La propuesta de Turing se basaba en un entretenimiento popular llamado «el juego de la imitación». En el juego original, un hombre y una mujer se ocultan a la vista y una tercera persona, el interrogador, les va haciendo preguntas alternativamente. A través de la lectura de sus respuestas tiene que determinar quién es el hombre y quién la mujer. La gracia está en que el hombre tiene que intentar engañar al interrogador, mientras que la mujer se esfuerza por ayudarlo, lo cual conduce a que las respuestas de uno y otro

lado resulten sospechosas. Turing reemplazó al hombre y a la mujer por un ordenador y una persona. Así nació lo que hoy conocemos como el «test de Turing»: un ordenador y una persona reciben las preguntas mecanografiadas de un juez humano, y si ese juez no logra identificar debidamente quién es el ordenador, el ordenador gana. Turing argumentó que, a partir de ese resultado, no dispondremos de ningún buen motivo para afirmar que la máquina carezca de inteligencia, sin importar que esta sea humana o no. Así, la cuestión de que la máquina disponga de inteligencia reemplazó la cuestión sobre si la máquina puede pensar de verdad.

El test de Turing, en realidad, es muy difícil: ningún ordenador lo ha superado. Por supuesto, en 1950 Turing desconocía este resultado a largo plazo; no obstante, al reemplazar las preguntas filosóficas problemáticas sobre la «consciencia» y el «pensamiento» con un test de resultados observables alentó la visión de la IA como una ciencia legítima con un objetivo bien definido. Mientras la IA cobraba forma durante los años cincuenta, muchos de sus pioneros y seguidores coincidieron con Turing: todo ordenador que pudiera mantener una conversación sostenida y convincente con una persona estaría, tal y como reconoceríamos la mayoría de nosotros, haciendo algo para lo que es necesario el pensamiento (sea eso lo que sea).

LA INTUICIÓN DE TURING / EL INGENIO COMO DISTINCIÓN

Turing se había labrado una reputación como matemático mucho antes de comenzar a escribir sobre IA. En 1936 publicó un artículo corto sobre el significado concreto de la palabra «computador», que en aquel momento se refería a la persona que seguía una serie de pasos para obtener un resultado definido (como la realización de un cálculo).² En aquel artículo reemplazó al computador humano por la idea de una máquina que realizara el mismo trabajo. El texto se adentraba en unas matemáticas de gran dificultad. Pero, mientras se refería a las máquinas, no hacía ninguna referencia al pensamiento humano ni a la mente. Las máquinas pueden operar de manera automática, afirmaba Turing, y los problemas que solucionan no requieren

de ninguna ayuda «externa» o inteligencia. Esa inteligencia externa —el factor humano— es lo que los matemáticos a veces denominan «intuición».

El trabajo que Turing dedicó en 1936 a las máquinas computadoras ayudó a lanzar la ciencia informática como disciplina, y representó una contribución importante a la lógica matemática. Aun así, al parecer Turing pensó que aquella definición temprana pasaba por alto una cuestión esencial. De hecho, la misma idea de que la mente o las facultades humanas pudieran ayudar a solucionar problemas apareció dos años después en su tesis doctoral, un inteligente pero fallido intento de esquivar uno de los resultados obtenidos por Kurt Gödel, matemático de origen austriaco especializado en lógica (volveremos a él en un rato). La tesis de Turing contiene este curioso pasaje sobre la intuición, que compara con otra capacidad mental a la que llama «ingenio»:

El razonamiento matemático puede ser considerado, de manera bastante esquemática, como un ejercicio de combinación entre dos facultades, a las que podríamos denominar intuición e ingenio. La actividad de la intuición consiste en realizar juicios espontáneos que no son el resultado de un hilo de razonamientos conscientes. Esos juicios son a menudo correctos, pero de ninguna manera lo son siempre (dejando de lado la cuestión sobre lo que se quiera decir con «correcto»). A menudo resulta posible encontrar otra manera de verificar la corrección de un juicio intuitivo. Por ejemplo, se puede juzgar que todos los números enteros positivos son factorizables en números primos; la argumentación matemática detallada conducirá a idéntico resultado. Esta también incluirá juicios intuitivos, pero serán menos susceptibles a la crítica que el juicio original sobre la factorización. No pretendo explicar esta idea de «intuición» de manera más explícita».

A continuación, Turing pasa a explicar el ingenio:

En matemáticas, el ejercicio del ingenio consiste en apoyar la intuición a través de una disposición adecuada de las proposiciones, y quizá de las figuras geométricas o de los dibujos. Lo que se pretende es que, cuando estos se encuentren dispuestos de manera verdaderamente correcta, la validez de los pasos intuitivos que sean necesarios no pueda ser motivo de una duda seria.³

Aunque su lenguaje se dirija a los especialistas, Turing señala lo evidente: por lo general, los matemáticos escogen sus problemas o «ven» un problema de interés en el que trabajar sirviéndose de una habilidad que cuando menos parece no poder dividirse en pasos, y que, por tanto, no se presta con claridad a la programación informática.

LA PERCEPCIÓN DE GÖDEL

También Gödel pensaba en la inteligencia mecánica. Igual que Turing, estaba obsesionado con la diferencia entre «ingenio» (mecánica) e «intuición» (mente). La distinción que él realizaba era en esencia la misma que la de Turing, aunque con un lenguaje diferente: demostración frente a verdad (o «teoría de la demostración» frente a «teoría de los modelos», en la jerga matemática). ¿Son, en definitiva, los de demostración y verdad el mismo concepto?, se preguntó Gödel. En caso afirmativo, las matemáticas e incluso la ciencia misma podrían entenderse de manera exclusivamente mecánica. Según esa visión, el pensamiento humano también sería mecánico. El concepto de AI, aunque el término no se hubiera acuñado aún, flotaba sobre la cuestión. ¿Se puede reducir la intuición de la mente, su capacidad para captar la verdad y el significado, a una máquina, a la computación?

Esta era la pregunta de Gödel. Al intentar contestarla, se encontró con un obstáculo que no tardaría en darle fama mundial. En 1931, Gödel publicó dos teoremas de lógica matemática conocidos como los teoremas de incompletitud. En ellos demostró las limitaciones inherentes a todos los sistemas matemáticos formales. Fue un golpe brillante. Gödel demostró de manera inconfundible que las matemáticas —toda la matemática, con ciertas suposiciones directas— no son, hablando en sentido estricto, ni mecánicas ni «formalizables». De manera más específica, Gödel demostró que en todo sistema formal (matemático o informático) han de existir proposiciones Verdaderas, con uve mayúscula, pero que no se pueden comprobar dentro del sistema mismo, sirviéndose de alguna de sus normas. La mente humana puede reconocer esa proposición Verdadera, pero el sistema en el que se ha formulado no la puede demostrar (cosa que sí es demostrable).

¿Cómo alcanzó Gödel esa conclusión? Los detalles son técnicos y complicados, pero la idea básica de Gödel es que podemos tratar un sistema matemático lo bastante complejo para realizar sumas igual que un sistema de significado, casi como si fuera una lengua natural del estilo del inglés o el francés —y lo mismo sirve para todos los sistemas de mayor complejidad—. Al tratarlo de esa manera, posibilitamos que el sistema hable sobre sí

mismo. Y puede contarnos, por ejemplo, que presenta ciertas limitaciones. Esa fue la percepción de Gödel.

Los sistemas formales, como los que aparecen en las matemáticas, permiten la expresión precisa de verdades y falsedades. Por lo general, establecemos lo que es verdad utilizando las herramientas de la demostración —nos servimos de unas reglas para demostrar algo y así sabemos que es indudablemente cierto. Pero ¿hay proposiciones verdaderas que no se puedan demostrar? ¿Puede la mente saber cosas que se le escapen al sistema? En el sencillo caso de la aritmética, expresamos verdades escribiendo ecuaciones como « $2 + 2 = 4$ ». Las ecuaciones básicas son proposiciones verdaderas dentro del sistema aritmético, demostrables según las normas de la aritmética. Aquí se da una equivalencia entre lo demostrable y lo verdadero. Antes de Gödel, los matemáticos pensaban que la matemática entera presentaba esa propiedad. Eso implicaba que las máquinas podrían producir en serie todas las verdades de los diferentes sistemas matemáticos limitándose a aplicar las normas de manera correcta. Es una idea hermosa, pero no es cierta.

A Gödel se le ocurrió la extraña pero poderosa propiedad de la autorreferencia. Se puede formar una versión matemática de expresiones autorreferenciales como «Esta proposición no se puede demostrar dentro de este sistema» sin quebrantar las reglas de los sistemas matemáticos. Pero las denominadas «proposiciones autorreferenciales de Gödel» introducen contradicciones en la matemática: si son ciertas, son indemostrables. Si son falsas, puesto que afirman ser indemostrables, en realidad son ciertas. Lo verdadero significa falso, y lo falso, verdadero: es una contradicción.

Retomando el concepto de «intuición», nosotros, los seres humanos, podemos ver que, de hecho, la proposición de Gödel es verdadera, pero, por culpa del resultado de Gödel, sabemos también que las normas del sistema no pueden demostrarla —en efecto, el sistema se muestra ciego ante aquello que sus normas no alcanzan a cubrir—. ⁴ Lo que es verdad y lo que es demostrable se desmontan entre sí. Y es posible que pase lo mismo con la mente y la máquina. En cualquier caso, los sistemas puramente formales tienen sus límites. No pueden probar desde su propio lenguaje algo que es cierto. En otras palabras, nosotros podemos ver cosas que al ordenador se le escapan. ⁵

El resultado de Gödel representó un duro golpe para la idea, popular en aquel momento, de que todas las matemáticas se podían convertir en

operaciones basadas en normas que produjeran una verdad matemática tras otra. El *Zeitgeist* pertenecía al formalismo, no a la conversación sobre las mentes, los espíritus, las almas y demás. En el campo de las matemáticas, el movimiento formalista señaló un giro más amplio de los intelectuales hacia el materialismo científico y, en particular, el positivismo lógico —un movimiento dedicado a erradicar la metafísica tradicional, como el platonismo, con esas formas abstractas que no se podían percibir con los sentidos, y las nociones tradicionales de la religión, como la existencia de Dios—. En efecto, el mundo estaba orientándose hacia la idea de las máquinas de precisión. Y nadie abrazó la causa formalista con tanto vigor como el matemático alemán David Hilbert.

EL DESAFÍO DE HILBERT

A principios del siglo xx (antes de Gödel), David Hilbert había lanzado un desafío al mundo matemático: demostrar que la totalidad de las matemáticas descansaba sobre un fundamento seguro. La ansiedad de Hilbert era comprensible. Si las normas puramente formales de la matemática no podían demostrar todas y cada una de sus verdades, al menos en teoría era posible que las matemáticas escondieran contradicciones y paparruchas. Que hubiera una contradicción oculta en algún lugar de las matemáticas lo arruinaba todo, porque a partir de una contradicción se puede demostrar cualquier cosa. Y, por tanto, el formalismo ya no servía para nada.

Hilbert expresó el sueño de cualquier formalista: demostrar al fin que las matemáticas eran un sistema cerrado y regido solo por normas. La verdad era tan solo una «demostración». Adquirimos conocimiento cuando nos limitamos a rastrear el «código» de una demostración y confirmamos que no se ha violado ninguna norma. El sueño más amplio de Hilbert, apenas disfrazado, apuntaba en realidad a una cosmovisión, a una imagen del universo en que este mismo fuera un mecanismo. La IA comenzó a cobrar forma como idea, una postura filosófica que también podía demostrarse. El formalismo trataba la inteligencia como si fuera un proceso reglado. Una máquina.

Hilbert lanzó su desafío durante el Segundo Congreso Internacional de Matemáticos, que se celebró en París en 1900. El mundo intelectual le dedicó su atención. El desafío constaba de tres partes principales: demostrar que las matemáticas eran una disciplina completa; demostrar que las matemáticas eran una disciplina consistente y demostrar que las matemáticas eran una disciplina decidible.

Con la publicación de sus teoremas de incompletitud, en 1931, Gödel hirió de muerte las partes primera y segunda del desafío de Hilbert. La cuestión de la decidibilidad quedó sin respuesta. Un sistema es decidible cuando existe un procedimiento definido (una demostración o una secuencia de pasos deterministas y evidentes) para establecer si una proposición construida a partir de las normas de ese sistema es verdadera o falsa. La proposición $2 + 2 = 4$ tiene que ser Verdadera, y la proposición $2 + 2 = 5$ tiene que ser Falsa. Y sucede lo mismo con todas las proposiciones que se puedan realizar con validez utilizando los símbolos y las reglas del sistema. Puesto que se creía que la aritmética era la base de las matemáticas, demostrar que las matemáticas eran decidibles implicaba demostrar el resultado de la aritmética y sus extensiones. Eso equivaldría a decir que los matemáticos, al «jugar» su partida con reglas y símbolos (la idea formalista), participaban de hecho en un juego válido que nunca conduciría a la contradicción ni al absurdo.

Turing quedó fascinado por el resultado de Gödel, que demostraba no el poder de los sistemas formales, sino más bien sus limitaciones. Se puso a trabajar en la parte que quedaba del desafío de Hilbert y comenzó a pensar en serio si podía existir un proceso de decisión para los sistemas formales. En 1936, con un artículo titulado «Números computables», demostró que no era así. Turing se dio cuenta de que el uso de la autorreferencia por parte de Gödel también podía aplicarse a las preguntas sobre los procesos de decisión, o, en efecto, a los programas informáticos. En especial, se percató de que debían de existir números (reales) que ningún método definido pudiera «calcular» al escribir su expansión decimal, dígito a dígito. Importó un resultado del matemático del siglo XIX Georg Cantor, quien había demostrado que los números reales (aquellos con expansión decimal) eran más numerosos que los enteros, por más que tanto los números reales como los enteros fueran infinitos. Es posible que Turing se subiera sobre hombros de gigantes, pero, al final, su labor en «Números computables» demostró una imposibilidad. Fue un resultado restrictivo: no era posible ningún

proceso de decisión universal. En otras palabras, las reglas —incluso en matemáticas— no bastan. Hilbert se había equivocado.⁶

LO QUE IMPLICÓ PARA LA IA

Lo importante de cara a la IA es lo siguiente: Turing refutó que las matemáticas fueran decidibles inventando una máquina, una máquina determinista, que no requería de ninguna intuición o inteligencia para resolver problemas. Hoy en día nos referimos a esa formulación abstracta de una máquina como la máquina de Turing. Ahora mismo estoy tecleando en una de ellas. Las máquinas de Turing son los ordenadores. Que el marco teórico de la informática se implementara como idea colateral, como un medio para obtener un fin diferente, es una de las grandes ironías de la historia intelectual. Mientras trabajaba para refutar que las matemáticas mismas fueran decidibles, Turing fue el primero en inventar algo preciso y mecánico: el ordenador.

En su tesis de 1938, Turing esperaba que los sistemas formales fueran ampliables incluyendo normas adicionales (y a continuación conjuntos de normas, y conjuntos de conjuntos de conjuntos de normas) que pudieran resolver el «problema de Gödel». Descubrió, en cambio, que aquel sistema nuevo y más potente tendría un problema de Gödel nuevo y más complejo. No había manera de sortear la incompletitud de Gödel. No obstante, enterrada bajo las complejidades del razonamiento de Turing sobre los sistemas formales, había una extraña sugerencia que resultaba relevante de cara a la posible existencia de la IA. ¿Y si la facultad de la intuición no se podía reducir a un algoritmo, a las normas de un sistema?

En su tesis de 1938, Turing intentaba encontrar una salida para el resultado restrictivo de Gödel, pero descubrió que era imposible. En su lugar cambió de marcha, se puso a explorar la manera en que, en sus propias palabras, podría «reducir en gran medida» el requisito de la intuición humana a la hora de realizar cálculos. Su tesis tomó en consideración el poder del ingenio al crear sistemas de normas cada vez más complicados (resultó que el ingenio podía volverse universal: hay máquinas capaces de tomar como referencia a otras máquinas y así dirigir todas las que se puedan construir. Esta percepción, técnicamente una

máquina de Turing universal en vez de una simple, iba a convertirse en el ordenador digital). Pero, en su trabajo formal sobre la computación, Turing se había ido de la lengua (quizá de manera involuntaria). Al permitir que la intuición fuera diferente y externa respecto a las operaciones de un sistema puramente formal como es el ordenador, Turing estaba, de hecho, sugiriendo que podían existir diferencias entre los programas de ordenador dedicados a las matemáticas y los matemáticos.

Por tanto, fue curioso el giro que Turing realizó entre sus primeros trabajos de los años treinta y la especulación de amplio espectro acerca de la posible aparición de ordenadores inteligentes en «Maquinaria computacional e inteligencia», que se publicó una década larga después. Hacia 1950, el debate sobre la intuición había desaparecido de los textos de Turing sobre las implicaciones de Gödel. Su interés se trasladó, en efecto, a la posibilidad de que los mismos ordenadores se convirtieran en «máquinas intuitivas». En esencia, decidió que el resultado de Gödel no era aplicable al asunto de la IA: si los seres humanos somos ordenadores muy avanzados, el resultado de Gödel solo implica que hay algunas proposiciones que no podemos comprender o ver como verdaderas, tal y como sucede con otros ordenadores menos complejos. Esas proposiciones podrían ser complejas e interesantes a extremos fantásticos. O tal vez fueran banales pero abrumadoramente complejas. El resultado de Gödel dejaba abierta la cuestión de si la mente no era más que una máquina de gran complejidad, con unas limitaciones muy complejas.

En otras palabras, la intuición había pasado a formar parte de las ideas de Turing acerca de las máquinas y sus poderes. El resultado de Gödel no podía afirmar (según Turing, en cualquier caso) que la mente fuera una máquina o no. Por un lado, la incompletitud sostiene que algunas proposiciones pueden entenderse como verdaderas desde el uso de la intuición, pero que eso no se puede demostrar a partir de un ordenador que se sirva del ingenio. Por el otro, un ordenador más poderoso puede utilizar axiomas (o más bits de código relevante) y demostrar el resultado, mostrando así que la intuición no está lejos de la computación en lo que a este problema se refiere. La cosa se convierte en una carrera armamentística: un ingenio cada vez más poderoso que sustituya a la intuición en problemas cada vez más complejos. Nadie puede anticipar quién ganará la carrera, así que nadie puede argumentar nada —usando el resultado de la incompletitud— sobre las diferencias inherentes entre

intuición (la mente) e ingenio (la máquina). Pero, tal y como Turing sin duda sabía, de ser eso cierto, también lo sería al menos la posibilidad de una inteligencia artificial.

Así, entre 1938 y 1950, Turing cambió de opinión acerca del ingenio y la intuición. En 1938, la intuición era el misterioso «poder de selección» que ayudaba a los matemáticos a decidir los sistemas con los que debían trabajar y los problemas que debían resolver. La intuición no era algo que se encontrara en el ordenador. Era algo que tomaba decisiones acerca del ordenador. En 1938, Turing no creía que la intuición formara parte de sistema alguno, lo cual sugería no solo que la mente y la máquina eran fundamentalmente diferentes, sino que una IA paralela al pensamiento humano resultaba casi imposible.

Sin embargo, para 1950 había cambiado de parecer. Con el test de Turing, desafió a los expertos e hizo una especie de defensa de la intuición en las máquinas; fue como si preguntara: «¿Por qué no?». Aquello supuso un cambio radical. Parecía que una nueva visión de la inteligencia comenzaba a cobrar forma.

¿Por qué ese cambio? Entre 1938 y 1950, a Turing le pasó algo ajeno al ámbito de las matemáticas estrictas y la lógica y los sistemas formales. Fue algo que le pasó, de hecho, a toda Gran Bretaña, y ciertamente a la mayor parte del mundo. Lo que pasó fue la segunda guerra mundial.

Capítulo 2

Turing en Bletchley

A Turing le fascinaba el juego del ajedrez —igual que a I. J. «Jack» Good, su colega matemático en tiempos de guerra. Cuando se enfrentaban (solía ganar Good), elaboraban procesos de decisión y reglas de oro para los movimientos ganadores. Jugar al ajedrez implica seguir las reglas del juego (ingenio), pero también parece requerir de cierta percepción (intuición) sobre las jugadas que pueden elegirse según las diferentes posiciones que se den sobre el tablero. Para ganar al ajedrez no basta con aplicar las reglas; en primer lugar, hay que saber qué reglas escoger.

Turing veía el ajedrez como una manera útil (y sin duda entretenida) de pensar sobre las máquinas y la posibilidad de conferirles intuición. Al otro lado del Atlántico, el fundador de la teoría de la información moderna, Claude Shannon, colega y amigo de Turing en Bell Labs, también pensaba en el ajedrez. Más adelante construyó uno de los primeros ordenadores que lo jugaron, una ampliación de la labor que había realizado anteriormente en un protoordenador llamado «el analizador diferencial», que podía convertir ciertos problemas de cálculo en procedimientos mecánicos.¹

EL PRINCIPIO DE LA SIMPLIFICACIÓN DE LA
INTELIGENCIA

El ajedrez fascinaba a Turing y a sus colegas en parte porque parecía que un ordenador podría programarse para jugar sin que la persona que lo programara necesitara saber todo por anticipado. Puesto que los dispositivos informáticos implementaban conectores lógicos como *si-entonces*, *o* e *y*, se podría ejecutar un programa (un conjunto de instrucciones) que generara resultados diferentes dependiendo de los escenarios con los que se encontrara mientras repasaba sus instrucciones. Esa capacidad para cambiar de rumbo según lo que «viera» parecía, a juicio de Turing y sus colegas, simular un aspecto fundamental del pensamiento humano.²

Los jugadores de ajedrez —Turing, Good, Shannon y demás— tenían también en la cabeza otro problema matemático con una apuesta mucho más elevada. Trabajaban para sus gobiernos, ayudando a descifrar los códigos secretos que usaba Alemania para coordinar sus ataques contra los barcos comerciales y militares que cruzaban el canal de la Mancha y el océano Atlántico. Turing se comprometió con un esfuerzo desesperado por ayudar a derrotar a la Alemania nazi durante la segunda guerra mundial, y fueron sus ideas sobre computación las que contribuyeron a alterar el curso de la guerra.

BLETCHLEY PARK

Bletchley Park, sita de manera discreta en un pueblo pequeño y alejado del reguero de bombas que caían sobre Londres y la Gran Bretaña metropolitana, era un centro de investigación establecido para ayudar a descubrir la localización de los *U-boote*, los submarinos alemanes, que causaban estragos en las rutas marinas del canal de la Mancha. Los submarinos nazis representaban un problema capital para las fuerzas aliadas; habían hundido miles de embarcaciones y destruido enormes cantidades de suministros y equipamiento. Para mantener el esfuerzo de guerra, Gran Bretaña necesitaba importaciones de treinta millones de toneladas al año. En un momento dado, los *U-boote* llegaron a reducir esa cantidad en 200.000 toneladas al mes, siguiendo una estrategia de guerra reveladora y potencialmente catastrófica, para la que durante bastante tiempo no hubo réplica. En respuesta, el gobierno británico reunió a un

grupo de criptoanalistas, jugadores de ajedrez y matemáticos talentosos para que investigaran la manera de descifrar las comunicaciones con los submarinos, conocidas como «cifrados». (Un «cifrado» es un mensaje oculto. Descifrar un mensaje consiste en convertirlo de nuevo en un texto legible.)³

Los códigos se generaban a través de un aparato con aspecto de máquina de escribir conocido como Enigma, que se comercializaba desde los años 1920 pero que los alemanes habían reforzado de manera importante para usarla en la guerra. Las máquinas Enigma modificadas se utilizaron en todo tipo de comunicaciones estratégicas dentro del esfuerzo de guerra nazi. La Luftwaffe, por ejemplo, las usó en su gestión de la guerra aérea, y lo mismo hizo la Kriegsmarine en sus operaciones navales. En general, se consideraba que los mensajes encriptados con la máquina Enigma modificada eran indescifrables.

El papel que Turing desempeñó en Bletchley y su consiguiente ascenso a la categoría de héroe nacional después de la guerra es una historia que ya se ha contado muchas veces. (En 2014, una gran producción cinematográfica, *The Imitation Game [Descifrando Enigma]*, dramatizó su trabajo en Bletchley, así como su rol consiguiente en el desarrollo de los ordenadores.) El mayor logro de Turing fue relativamente desaborido, según criterios matemáticos puros, porque explotó una vieja idea de la lógica deductiva. El método, al que él y otras personas se referían medio en broma como «turinguismo», se basó en eliminar amplios números de posibles soluciones para los códigos de Enigma encontrando combinaciones en las que hubiera contradicciones. Las combinaciones contradictorias son una imposibilidad; en un sistema lógico no puede darse «A» y «no A» a la vez, tal y como no podemos estar «en la tienda» y «en casa» al mismo tiempo. El turinguismo fue una idea ganadora, y se convirtió en un gran éxito en Bletchley. Logró lo que se había exigido a aquellos «jóvenes genios» recluidos en el laboratorio de ideas al acelerar el descifrado de los mensajes de Enigma. Otros científicos de Bletchley concibieron estrategias diferentes para descifrar los códigos.⁴ Sus ideas se ponían a prueba con una máquina llamada Bombe —nombre burlón que provenía de una máquina polaca anterior, la Bomba, y que con toda probabilidad se inspiró en los ruiditos que esta realizaba al terminar cada uno de sus cálculos—. Pensemos en la Bombe como en un protoordenador, capaz de ejecutar diferentes programas.

Más o menos en 1943, el Eje perdió su ventaja bélica en beneficio de las fuerzas aliadas, y ello se debió en no poca medida al esfuerzo continuado de los descifradores de Bletchley. Aquel equipo obtuvo un éxito célebre, y sus miembros se convirtieron en héroes de guerra. Hicieron carrera. Bletchley, mientras tanto, también se reveló como un refugio para el pensamiento dedicado a la computación: la Bombe era una máquina que ejecutaba programas para resolver problemas que los seres humanos por sí mismos no podían solucionar.

¿MÁQUINAS INTUITIVAS? NO

En el caso de Turing, Bletchley desempeñó un papel capital de cara a que materializara sus ideas sobre la posibilidad de crear máquinas inteligentes. Igual que sus colegas Jack Good y Claude Shannon, Turing percibió el poder y la utilidad de sus «juegos mentales» como criptoanalistas durante la guerra: podían descifrar mensajes que de otro modo resultaban completamente opacos para los militares. Los nuevos métodos computacionales no solo resultaban interesantes para pensar en un juego de ajedrez automatizado, sino que podían, de manera bastante literal, hundir barcos de guerra.

Turing (una vez más) pensaba en una abstracción: mentes y máquinas, o la idea general de inteligencia. Pero había algo extraño en su visión de lo que aquello implicaba. En los años cuarenta, la inteligencia no era un rasgo que se atribuyera en general a los sistemas formales, como era el caso de la Bombe de Bletchley, una máquina descifradora puramente mecánica. Gödel había demostrado que, por norma, la verdad no podía reducirse a lo formal, en el sentido de que participara en un juego formal con un conjunto de reglas establecidas, pero recuerda que su demostración dejó abierta la cuestión de si una máquina específica podría incorporar la intuición de la que se sirve la mente para tomar decisiones sobre las reglas que se deben seguir, pese a que no pudiera existir ningún sistema supremo capaz de demostrarlo todo (tal y como el propio Gödel había revelado de manera tan definitiva en 1931).

Tras abandonar Bletchley, Turing dedicó cada vez más tiempo a la cuestión de si era posible construir una máquina que fuera lo bastante

potente como para usar a la vez la intuición y el ingenio. El enorme número de combinaciones posibles que había que comprobar de cara a descifrar los códigos alemanes resultaba abrumador para la intuición humana. Pero unos sistemas que contaran con los programas adecuados podrían cumplir con esa tarea al simplificar aquellas vastas posibilidades matemáticas. Para Turing, eso sugería que la intuición podía cobrar cuerpo en las máquinas. En otras palabras, el éxito de Bletchley implicaba que quizá se pudiera construir una inteligencia artificial.

No obstante, para que esa línea de pensamiento cobrara sentido, había que decidirse por una idea concreta de «inteligencia». La inteligencia, tal y como la ejercen los seres humanos, debía ser reducible —analizable— según los términos de la capacidad de la máquina. En esencia, la inteligencia debía ser reducible a la forma de la resolución de problemas. Al fin y al cabo, en eso consiste el juego del ajedrez y en eso consiste también descifrar un código.

Y ahí está: la mayor muestra de genio por parte de Turing, y también su mayor error, consistió en pensar que la inteligencia humana se limitaba a resolver problemas. Tanto si había explicitado las ideas sobre máquinas inteligentes de su «Maquinaria computacional e inteligencia» de 1950 durante los años de la guerra como si no, queda claro que la experiencia de Bletchley materializó su visión posterior sobre la IA, y queda claro que la IA, a su vez, siguió de cerca esa misma senda, aunque sin el autoanálisis que hubiera sido necesario.

Pero una mirada más atenta al éxito descodificador de Bletchley revela de manera inmediata una simplificación peligrosa en sus ideas filosóficas acerca del hombre y la máquina. Bletchley fue un sistema inteligente, resultado de la coordinación militar (incluyendo el espionaje y la inteligencia, así como la captura de naves enemigas); de la inteligencia social que se estableció entre los militares y los diversos científicos e ingenieros que había allí, y (como sucede con todo en esta vida) a veces fue también cuestión de pura suerte. Lo cierto es que, en cuanto realidad práctica, la máquina Enigma modificada por los alemanes era impenetrable por medios puramente mecánicos. Los alemanes eran conscientes de ello; se habían basado en argumentos matemáticos sobre las dificultades de la descodificación mecánica. En parte, el éxito de Bletchley se debió, irónicamente, a la tozuda confianza de los comandantes nazis en el carácter inexpugnable de los cifrados de Enigma —de modo que, en momentos

cruciales, tras descubrir que ciertos mensajes habían sido descifrados, se negaron a modificar o reforzar las máquinas, echándoles la culpa a operaciones de espionaje encubiertas en vez de aceptar aquella derrota científica—. Pero la niebla de la guerra hace que se mezclen no solo diferentes tecnologías novedosas, sino nuevas formas de inteligencia humana y social. La guerra no es como el ajedrez.

Por ejemplo, al principio de la guerra, las fuerzas polacas recuperaron fragmentos importantes de comunicaciones de Enigma que más tarde revelaron pistas de valor incalculable para la labor de Bletchley. Los polacos habían usado esos fragmentos (junto con otros, procedentes de fuentes rusas) para desarrollar su propia Bombe, aunque más simple, en una fecha tan temprana como 1938. La versión muy mejorada de Turing a principios de 1940 —la Bombe que usaba su «turinguismo»— dependió de aquella primeriza labor polaca, facilitada por los hechos que tenían lugar en el campo de batalla. Turing también vio a su colega Gordon Welchman introducir cambios en su propio diseño, al que le añadió un «tablero diagonal» para simplificar aún más la búsqueda de contradicciones,⁵ como respuesta a las mejoras que los alemanes habían realizado en Enigma. Ahí había dos mentes humanas sirviéndose de la intuición, trabajando conjuntamente en sociedad.

Hubo otros acontecimientos en el teatro de la guerra que resultaron de importancia capital. El 8 de junio de 1940, un portaaviones británico se hundió delante de la costa noruega. Aquel ataque facilitó la localización de los *U-Boote*, si bien se cobró un precio elevado con los numerosos marineros que acabaron en el fondo del mar. Pocas semanas antes, a finales de abril de 1940, la patrullera alemana VP2623, un miembro de la flota especialmente devastador, había sido capturada con un tesoro de pruebas de Enigma en su interior. Las piezas que se necesitaban para resolver el puzle de Enigma estaban llegando a manos aliadas y se abrían camino hacia el grupo de Bletchley.

Por sí mismos, aquellos fragmentos resultaban por completo inadecuados para descifrar con rapidez el futuro de las comunicaciones alemanas; para los criptoanalistas de Bletchley no eran más que «conjeturas», según la definición de un biógrafo de Turing. Pero facilitaron un primer paso de radical importancia a la hora de dar con la manera de programar las máquinas Bombe. Turing y sus colegas lo denominaron «la ponderación de las pruebas», tomando prestado un término que acuñó el científico y lógico

norteamericano C. S. Peirce (quien ocupará un papel destacado en la segunda parte de este libro).⁶

Los matemáticos interpretan el peso de la evidencia de maneras diferentes, pero, en el caso del éxito de Bletchley (y para asuntos más amplios relacionados con la IA), equivale a aplicar conjeturas informadas, o intuiciones, para dirigir el ingenio, o las máquinas. Un fragmento de texto descifrado procedente de un submarino capturado puede significar cualquier cosa, tal y como una bola blanca hallada cerca de una bolsa de bolas blancas puede significar cualquier cosa, pero en cada caso podemos realizar suposiciones inteligentes para comprender lo que ha sucedido. Pensamos que resulta muy probable que la bola blanca haya salido de esa bolsa, pese a no haber visto que la sacaran de ella. Se trata de una suposición. No se puede demostrar que ese tipo de suposiciones sean ciertas, pero, cuanto mejor funcione la intuición humana al establecer las condiciones iniciales para trazar los procesos mecánicos, mejores serán las posibilidades de que esos procesos acaben obteniendo los resultados deseados en vez de, pongamos por caso, prolongarse sin rumbo fijo, siguiendo direcciones erróneas o engañosas. El peso de la evidencia —suponer— hizo que las Bombe funcionaran.

Los científicos de Bletchley no se limitaron a proveer a las Bombe de información, dejándolas luego para que realizaran la labor incansable e importante de eliminar millones de códigos o cifrados incorrectos. Desde luego, las Bombe fueron necesarias —eso es lo que Turing entendió con gran claridad, y lo que sin duda inundó su imaginación con la posibilidad de que aquellos «procesos mecánicos» pudieran reproducir o reemplazar a la inteligencia humana—, pero la realidad fue que el grupo de Bletchley se ocupó ante todo de hacer conjeturas. Al reconocer las pistas escondidas en el mosaico de instrucciones incompletas, cifrados y mensajes procedentes del campo de batalla, pasaron a generar hipótesis. En la ciencia, las conjeturas se definen así, como la formación de hipótesis (concepto que también utilizó Charles Sanders Peirce), y tienen una importancia fundamental para el progreso del saber humano. No es de extrañar, pues, que la obra de Bletchley equivaliera a un sistema de conjeturas acertadas. Su condición *sine qua non* no fue de tipo mecánico, sino que más bien podríamos describirla como una observación inicial inteligente. Las Bombe necesitaban que las apuntaran hacia algo, y que a continuación las impulsaran en ese sentido.

En sintonía con un tema que exploraremos en la segunda parte del libro, Peirce reconoció muy al principio, a finales del siglo XIX, que todas las observaciones que dan cuerpo a las ideas complejas y juicios de la inteligencia comienzan con una suposición, o lo que él llamó una abducción:

Al mirar por la ventana en esta hermosa mañana de primavera veo una azalea en plena floración. ¡No, no! No es eso lo que he visto, aunque sí se trate de la única manera en que puedo describirlo. Es una proposición, una frase, un dato; pero lo que percibo no es una proposición, una frase, un dato, sino apenas una imagen que yo hago inteligible en parte a través de la exposición de un hecho. Esa exposición es abstracta, pero lo que yo veo es concreto. Realizo una abducción cada vez que expreso cualquier cosa que haya visto en una frase. La verdad es que el entramado al completo de nuestro conocimiento es un fieltro opaco de hipótesis puras confirmadas y refinadas a través de la inducción. No se puede obtener el menor avance en el campo del conocimiento sin realizar una abducción a cada nuevo paso, o de otro modo nos quedaríamos mirando las cosas con expresión vacía.⁷

Turing y sus colegas de Bletchley comenzaron a ganar una guerra que había orbitado de los mandamases a los servicios de inteligencia gracias al uso, en efecto, de abducciones inteligentes en cada nuevo paso del camino. Hasta cierto punto, es evidente que Turing era consciente de ello (recordemos el debate sobre la intuición en la tesis que dedicó en 1938 a los números ordinales), pero no parece haber tenido un efecto apreciable en sus ideas posteriores sobre la naturaleza de la inteligencia y la posibilidad de crear máquinas inteligentes. Por brillante que se mostrara, formuló una simplificación de la inteligencia real. Se liberó del concepto que tanto le había subyugado con anterioridad: el de la intuición. El de las conjeturas.

SOBRE LA INTELIGENCIA SOCIAL (UNA ACOTACIÓN IMPORTANTE)

La inteligencia social también quedó visiblemente fuera de la forma en que Turing resolvió el acertijo sobre la inteligencia. Esto es de la mayor importancia de cara a comprender el desarrollo futuro de la IA. Por ejemplo, a Turing le desagradaba considerar que el pensamiento o la inteligencia pudieran ser circunstancias sociales o situacionales.⁸ Sin embargo, el éxito de Bletchley formó parte, en realidad, de un vasto sistema

que se extendió mucho más allá de las cuatro paredes del lugar. Se había puesto en marcha un esfuerzo inmenso, que no tardaría en atraer a Estados Unidos y la labor de científicos como Shannon, en Bell Labs, así como de los que trabajaron en el célebre Instituto de Estudios Avanzados de Princeton —donde tenían puestos Einstein, Gödel y John von Neumann—. El sistema expandido de máquinas humanas resulta en realidad mucho más realista como modelo de la manera en que se solucionan los problemas del mundo real —entre ellos, el de una guerra mundial debe contarse sin duda como uno de los más complejos e importantes.

La falta de oído musical de la IA para la inteligencia social o situacional ya se había comentado antes, y en tiempos más recientes lo ha hecho el científico especializado en aprendizaje automático François Chollet, quien lo resume bien en su crítica a la visión que tenía Turing sobre la inteligencia (y, de manera más amplia, a la del campo de la IA). Primero, la inteligencia es *situacional*, no existe nada parecido a una inteligencia general. Tu cerebro es una pieza dentro de un sistema más amplio que incluye tu cuerpo, tu entorno, a otros seres humanos y la cultura en su conjunto. Segundo, la inteligencia es *contextual*: lejos de existir en el vacío, cualquier inteligencia individual siempre se hallará definida a la vez que limitada por su entorno. (Y, en estos momentos, es el entorno, y no el cerebro, el que actúa como cuello de botella para la inteligencia.) Tercero, la inteligencia humana se encuentra en gran medida *externalizada*, contenida no en tu cerebro sino en tu civilización. Pensemos en los individuos como si fueran herramientas cuyos cerebros son módulos de un sistema cognitivo mucho más amplio que ellos mismos, un sistema que lleva mucho tiempo evolucionando.⁹

Según lo expresa Turing, la intuición se puede programar en una máquina, pero Chollet y críticos similares aseguran que esta no podrá alcanzar el nivel de la inteligencia humana. De hecho, la idea de programar la intuición ignora un aspecto fundamental de nuestros propios cerebros. Los seres humanos disponemos de inteligencia social. Disponemos de inteligencia emocional. Usamos nuestras mentes para algo más que para resolver problemas y acertijos, por complejos que sean (o, más bien, *sobre todo* cuando esos problemas son complejos).

La evidencia sugiere que Turing rechazó con firmeza esa visión de las personas, y en su lugar llegó a creer que la totalidad del pensamiento humano se podía entender, en efecto, desde el «desciframiento» de unos

«códigos» —o resolución de acertijos— y la práctica de juegos como el ajedrez. Lo importante es que, en algún momento de los años cuarenta, después de trabajar en Bletchley y sin duda durante la época en la que escribió el artículo aparecido en 1950 donde prefiguraba la IA, el pensamiento de Turing se decantó por una visión simplificada de la inteligencia. Fue un error atroz, que además se ha ido transmitiendo de generación en generación de científicos de IA hasta llegar al día de hoy.

EL ERROR DE TURING CON LA INTELIGENCIA Y UNA IA DÉBIL

Esa visión de la inteligencia como algo que resuelve problemas ayuda a explicar la producción de aplicaciones invariablemente débiles a lo largo de la historia de la IA. Los juegos, por ejemplo, han sido una fuente constante de inspiración para el desarrollo de técnicas avanzadas de IA, pero estos no dejan de ser versiones simplificadas de la vida que recompensan visiones también simplificadas de la inteligencia. Un programa de ajedrez puede desempeñarse bien en ese juego, pero se le dará bastante mal conducir un coche. El sistema Watson de IBM juega al *Jeopardy!*, pero no al ajedrez ni al go, y se requiere un esfuerzo inmenso de programación o de «conversión» para que la plataforma Watson realice otras funciones de extracción de datos y procesamiento del lenguaje natural, como con sus recientes (y en gran medida fallidas) incursiones en el terreno de la salud.

Por consiguiente, tratar la inteligencia como algo que resuelve problemas conduce a que las aplicaciones de la IA sean débiles. Sin duda, Turing fue consciente de ello, y en su artículo de 1950 especuló con la posibilidad de que se pudiera hacer que las máquinas aprendieran y así superar las limitaciones que surgen como consecuencia natural del diseño de unos sistemas informáticos que solo sirven para solucionar problemas. Si las máquinas aprendieran a volverse genéricas, seríamos testigos de una transición fluida entre las aplicaciones específicas y unos seres dotados de pensamiento general. Llegaríamos a la IA.

No obstante, el conocimiento que tenemos hoy choca con violencia contra el enfoque de aprendizaje sugerido de manera temprana por Turing.

Para alcanzar sus objetivos, los que en la actualidad denominamos «sistemas de aprendizaje automático» deben aprender algo específico. Los investigadores hablan de darle a la máquina un «sesgo» (sin las connotaciones negativas que le otorgamos en nuestra sociedad; no se pretende decir que la máquina sea cabezota o que cueste discutir con ella, ni que tenga motivaciones secretas, según el sentido habitual de la palabra). En el aprendizaje automático, el sesgo significa que el sistema está diseñado y puesto a punto para aprender algo. Pero, por supuesto, ese es precisamente el problema de producir aplicaciones débiles que resuelvan problemas. (Y es el motivo, por ejemplo, por el que los sistemas de aprendizaje profundo que usa Facebook para reconocer rostros humanos no han aprendido también a hacerte la declaración de la renta.)

Peor incluso, los investigadores se han dado cuenta de que darle a un sistema de aprendizaje automático un sesgo a la hora de aprender una tarea o aplicación concreta lleva a que tenga un rendimiento peor en otras tareas. Hay una correlación inversa entre el éxito de la máquina al aprender algo y que consiga aprender otra cosa. Incluso tareas en apariencia similares presentan esa relación inversa en su desempeño. Un sistema informático que aprenda a jugar al go a nivel de campeonato no aprenderá además a jugar al ajedrez a ese mismo nivel. El sistema del go ha sido diseñado de manera específica, con un sesgo particular hacia el aprendizaje de las reglas del go. Su curva de aprendizaje, tal y como la llaman, sigue, por tanto, el tanteo conocido de ese juego en particular y, en relación con cualquier otro juego, pongamos el *Jeopardy!* o el ajedrez, se vuelve inútil —de hecho, no existe.

El sesgo en el aprendizaje automático se ha entendido por lo general como una fuente de errores de aprendizaje, un problema técnico. (También, al ajustarse al uso común del lenguaje, ha adoptado acepciones secundarias que ofrecen resultados involuntarios pero inaceptables por, pongamos, su carga racial o de género.) El sesgo en el aprendizaje automático puede introducir errores solo porque el sistema no «busca» ciertas soluciones en primer lugar. Pero, de hecho, el sesgo es necesario para el aprendizaje automático: forma parte de él.

Un célebre teorema conocido como «no free lunch» demuestra con exactitud lo que observamos de manera anecdótica al diseñar y construir un sistema de aprendizaje. El teorema sostiene que, al aplicarse sobre un problema arbitrario, cualquier sistema de aprendizaje libre de sesgo no

obtendrá resultados mejores que los que proporciona el azar. Es una manera elegante de decir que los diseñadores de sistemas deben conferir a estos un sesgo de manera deliberada, para que aprendan su propósito. Tal y como señala el teorema, un sistema en verdad libre de sesgo no sirve para nada. Hay técnicas complicadas, como la del «preentrenamiento» con datos, que se sirven de métodos no supervisados que exponen los rasgos de los datos que hay que aprender. Todo ello forma parte integral de un aprendizaje automático exitoso. Lo que queda fuera del debate, no obstante, es que ajustar un sistema para que aprenda su propósito inculcándole el sesgo deseado implica que se vuelva restrictivo, en el sentido de que ya no podrá generalizarse a otros dominios. En parte, construir e implementar con éxito un sistema de aprendizaje automático lleva a que este no se encuentre libre de sesgo y no sea general, sino que se centre en un problema de aprendizaje particular. Visto así, la restricción se encuentra integrada hasta cierto punto en esos enfoques. El éxito y la restricción son las dos caras de una misma moneda.

Por sí solo, ese hecho ya arroja serias dudas sobre cualquier expectativa de progresión fluida entre la IA actual y la IA de nivel humano el día de mañana. La gente que asume que la ampliación de los métodos modernos de aprendizaje automático, como el aprendizaje profundo, podrán formarse desde cero o aprender a ser tan inteligentes como los seres humanos, no comprende las limitaciones fundamentales ya conocidas. Admitir la necesidad de suministrar un sesgo a los sistemas de aprendizaje es equivalente a la observación por parte de Turing de que la mente humana debe suministrar percepciones matemáticas externas a los métodos formales, ya que el sesgo del aprendizaje automático está determinado, antes del aprendizaje, por sus diseñadores humanos.¹⁰

EL LEGADO DE TURING

Para resumir la cuestión, la visión de la inteligencia como algo que resuelve problemas genera de manera necesaria aplicaciones débiles y, por tanto, resulta inadecuada para los objetivos más amplios de la IA. Heredamos esa visión de la inteligencia de Alan Turing. (¿A cuento de qué, por ejemplo, usamos el término «inteligencia artificial» en vez de, quizá, el de

«simulación de tareas humanas»?)¹¹ La genialidad de Turing consistió en deshacerse de los obstáculos y objeciones teóricas en el trayecto hacia la posibilidad de diseñar una máquina autónoma, pero con ello limitó el alcance y la definición de la inteligencia misma. No es de extrañar, pues, que la IA comenzara produciendo aplicaciones débiles de resolución de problemas y que haya seguido haciéndolo hasta el día de hoy.

Una vez más, a Turing le desagradaba la consideración del pensamiento o la inteligencia como algo social o situacional. Sin embargo, pese a su tendencia a entender la inteligencia humana como un proceso mecánico individual —lo que dio pie en los años cuarenta, con la aparición de los primeros ordenadores, a incontables menciones periodísticas al «cerebro mecánico»—, resulta evidente que la conversación acerca de la inteligencia implica siempre, y de manera necesaria, situarla en un contexto más amplio. La inteligencia general (no débil) del tipo que todos exhibimos a diario no se debe a ningún algoritmo que se esté ejecutando dentro de nuestras cabezas, sino que recurre a la totalidad del contexto cultural, histórico y social desde el que pensamos y actuamos en el mundo. La IA apenas habría avanzado si sus diseñadores hubieran abrazado un entendimiento de la inteligencia tan amplio y complejo —es cierto—. A la vez, a resultas de la simplificación realizada por Turing, hemos acabado usando aplicaciones débiles y no tenemos ningún motivo para esperar otras más generales si antes no se produce una reconceptualización radical de lo que queremos decir al hablar de IA.

En su artículo de 1950, Turing anticipó algunas de esas dificultades al sugerir que se podía hacer que las máquinas aprendieran. No obstante, lo que sabemos en este momento (en oposición a la excitación reciente sobre el aprendizaje automático) es que el aprendizaje mismo es un tipo de resolución de problemas posibilitado solo por la introducción de un sesgo en el aprendiz; sesgo que de manera simultánea facilita el aprendizaje de una aplicación en particular mientras que limita el desempeño en otras aplicaciones. De hecho, los sistemas de aprendizaje también son sistemas restrictivos de resolución de problemas. Puesto que no existe ningún puente teórico entre esos sistemas tan restrictivos y la inteligencia general del tipo que exhibimos los seres humanos, la IA ha caído en una trampa. Una serie de errores tempranos en la comprensión de la inteligencia han conducido, en grado diverso, pero de manera inexorable, a un punto muerto teórico en el núcleo de la IA.

Consideremos de nuevo la distinción original de Turing entre intuición e ingenio. Para él, el problema de la IA consistía en que la intuición — aquello que viene suministrado por el diseñador del sistema— pudiera de hecho «introducirse» en la parte formal de ese sistema (el ingenio de la máquina) y crear así un sistema capaz de escapar a la maldición de lo restrictivo al servirse de la intuición para escoger sus propios problemas — para volverse cada vez más inteligente y aprender—. Hasta el momento, nadie ha conseguido hacer eso con ningún ordenador. Nadie tiene la menor idea, siquiera, de la manera en que se podría llegar a ello. Sabemos que los diseñadores usan la intuición desde fuera de los sistemas de IA para indicar a estos los problemas específicos que deben resolver (o que deben aprender a resolver). La cuestión de que los sistemas utilicen la intuición de manera autónoma apunta directamente al núcleo de lo que denominaré «el problema de la inferencia», sobre el que hablaremos en la segunda parte del libro.

En esa segunda parte también habrá muchas otras cosas que comentar sobre «la trampa de la restricción». Pero antes hay más terreno que cubrir en este primer apartado. A continuación, pasaremos al tema de la superinteligencia, otro error de la inteligencia y una extensión natural del primero.

Capítulo 3

El error de la superinteligencia

Jack Good, el colega descifrador de Turing, también quedó fascinado por la idea de una máquina dotada de inteligencia. Es evidente que Turing allanó su imaginación cuando ambos se encontraban en Bletchley, pero, más adelante, Good añadió un giro como de ciencia ficción a las ideas de Turing acerca de la posibilidad de una inteligencia de nivel humano en los ordenadores. La idea de Good era sencilla: si una máquina puede alcanzar el nivel de la inteligencia humana, también puede sobrepasar el mero pensamiento humano.

Good consideraba evidente que una especie de bucle de retroalimentación permitiría a las máquinas inteligentes examinarse a sí mismas y mejorar, lo que conduciría a la creación de máquinas aún más inteligentes y resultaría en una «explosión de inteligencia» desenfrenada. La explosión de inteligencia seguiría porque cada máquina haría una copia aún más inteligente de sí misma, con el resultado de una curva exponencial de inteligencia en las máquinas que no tardaría en sobrepasar a los grandes genios de la humanidad. Good lo denominó «ultrainteligencia»:

Defínase como ultrainteligente a aquella máquina que puede rebasar con amplitud las actividades intelectuales de cualquier hombre, por listo que este sea. Puesto que el diseño de las máquinas es una de esas actividades intelectuales, la máquina ultrainteligente podría diseñar máquinas aún mejores; es incuestionable que seguiría una «explosión de inteligencia» y que la inteligencia humana quedaría muy atrás. Así, la primera máquina ultrainteligente sería el último producto que la humanidad necesitaría inventar, contando con que la máquina se mostrara lo bastante dócil como para indicarnos la manera de mantenerla bajo control.¹

El filósofo de Oxford Nick Bostrom iba a recuperar el tema de Good algunas décadas más tarde con un *best seller* titulado *Superinteligencia: Caminos, peligros, estrategias*, donde defendía la misma postura: que la consecución de la IA tendría como consecuencia el inicio de una inteligencia superior a la humana en un proceso cada vez más vertiginoso de automodificaciones. Con un lenguaje ominoso, Bostrom se hace eco del futurismo de Good acerca del advenimiento de las máquinas superinteligentes:

Ante la perspectiva de una explosión de inteligencia, nosotros, los seres humanos, somos como niños pequeños que juegan con una bomba. Esa es la disparidad entre la potencia de nuestro juguete y la inmadurez de nuestra conducta. La superinteligencia representa un desafío para el que no estamos preparados en este momento, y para el que no estaremos preparados durante mucho tiempo. No sabemos bien cuándo tendrá lugar la detonación, aunque si nos llevamos el artefacto a la oreja podemos oír un tictac amortiguado. En el caso de un niño que tuviera una bomba sin detonar entre las manos, lo más sensato sería que la dejara en el suelo con cuidado, que saliera rápidamente de la habitación y que llamara al adulto más cercano. Sin embargo, aquí nos encontramos no con un niño sino con muchos, y cada uno de ellos tiene acceso a un mecanismo de activación independiente. Las posibilidades de que todos demos el sentido común de deshacernos de algo tan peligroso parecen nimias. Siempre habrá algún idiota dispuesto a apretar el botón de ignición solo para ver qué pasa.²

Para Bostrom, la superinteligencia no es objeto de especulación ni una nebulosa, sino que se parece más a la llegada de las armas nucleares: un *hecho consumado* que tendrá consecuencias profundas y quizá nefastas para la raza humana. El mensaje está claro: no discutas si la superinteligencia va a llegar. Prepárate para su llegada.

¿Qué hemos de decir ante eso? El argumento de Good y Bostrom —la posibilidad de una máquina superinteligente— parece verosímil a primera vista. Pero, de manera poco sorprendente, nunca se especifica el mecanismo por el que la inteligencia de partida llevará a la superinteligencia. Good y Bostrom parecen tomarse la posibilidad de la superinteligencia como algo tan creíble y evidente que esta no requiere de mayores explicaciones. Pero sí que las necesita; tenemos que comprender el «cómo».

Si imaginamos una simple mejora como la de un *hardware* superior, la propuesta resulta demasiado trivial y ridícula como para que le dediquemos más tiempo. Es probable que ni siquiera un creyente incondicional en la inexorabilidad del progreso como Ray Kurzweil vaya a rebajar la inteligencia hasta ese punto —no pensamos que, al añadirle memoria RAM a un MacBook, estemos volviéndolo más inteligente (de verdad)—. El aparato irá más rápido, y podrá cargar aplicaciones de mayor tamaño y tal.

Pero, si por inteligencia entendemos algo interesante, esta tiene que hacer algo más complejo que cargar aplicaciones a gran velocidad. Esa parte más difícil de la inteligencia se queda sin comentar.

O supongamos que tomamos prestado el lenguaje del mundo de la biología (como la IA hace tan a menudo) y, a continuación, declaramos con seguridad que la capacidad computacional no involuciona, sino que evoluciona. Una mirada más profunda nos hará ver que ese argumento se ve afectado una vez más por una visión inadecuada e ingenua de la inteligencia. El problema —una omisión flagrante— es que no existe ninguna evidencia en el mundo biológico de que algún ser inteligente haya diseñado alguna vez una versión más inteligente de sí mismo. Los seres humanos somos inteligentes, pero a lo largo de la historia de la humanidad no hemos construido ninguna versión más inteligente de nosotros mismos.

Hay un requisito previo a la construcción de un cerebro más listo, y consiste en comprender el carácter cognitivo del que ya tenemos, en el sentido de que podemos imaginar escenarios, albergar pensamientos y sus conexiones, encontrar soluciones y descubrir nuevos problemas. Nos pasan cosas, razonamos a través de nuestras observaciones y de lo que ya sabemos, las respuestas brotan en nuestra cabeza. Todo ese zumbido de magia biológica sigue siendo opaco, la mayor parte de su «procesamiento» continúa pendiente de cartografiar. Y, sin embargo, llevamos milenios considerando e investigando nuestros procesos mentales y nuestras funciones cerebrales.

Siendo evidente que con nosotros no sucede así, ¿por qué debería una máquina mayormente inteligente desarrollar de golpe una percepción sobre sus propias capacidades cognitivas globales? Y, aunque lograra hacerlo, ¿cómo podría la máquina usar ese conocimiento para volverse más lista?

No es una cuestión de autosuperación. Podemos, por ejemplo, volvernos más inteligentes leyendo libros o yendo a la escuela; al formarnos posibilitamos un mayor desarrollo intelectual, etcétera. Todo eso no es motivo de controversia. Y nada de ello viene al caso. Un problema capital de las conjeturas acerca del aumento de inteligencia en los círculos de la IA es el de su carácter circular: hace falta una inteligencia (en apariencia general) para hacer crecer la inteligencia general. Si lo observamos con detenimiento veremos que no hay ninguna progresión lineal, solo misterio.

VON NEUMANN Y LAS MÁQUINAS AUTORREPLICANTES

Good introdujo a mediados de los años sesenta la idea de una IA evolutiva que condujera a la ultrainteligencia. Pero, casi dos décadas antes, John von Neumann ya había considerado esa idea y la había descartado. En una conferencia de 1948 en el Instituto de Estudios Avanzados de Princeton, Von Neumann explicó que, mientras que la reproducción humana a menudo mejora los «diseños» previos, resulta evidente que las máquinas con el cometido de diseñar máquinas nuevas y mejores se enfrentan a un escollo fundamental, ya que cualquier diseño para una máquina nueva tiene que aparecer especificado en la máquina madre. La máquina madre tendría que ser por necesidad más compleja que su creación, no menos: «La organización sintetizadora ha de ser por necesidad más compleja, de un orden superior, que la organización sintetizada», dijo.³

En otras palabras, Von Neumann señaló una diferencia fundamental entre la vida orgánica tal y como la conocemos y las máquinas que construimos. La predicción de la ultrainteligencia que había hecho Jack Good era un poco de ciencia ficción.

Von Neumann propuso que la máquina autorreplante debía tener, como mínimo, ocho partes, incluyendo un órgano que reciba y transmita «estímulos», un órgano «de fusión» que conecte las partes entre sí, un órgano «de corte» que interrumpa esas conexiones y un «músculo» para poder desplazarse. A continuación, bosquejó mecanismos que pudieran generar mejoras cognitivas de manera verosímil y que incluyeran un elemento aleatorio, similar a las mutaciones biológicas, para que permitiera las modificaciones necesarias. Pero Von Neumann pensó que, en vez de hacer avanzar el pensamiento de la máquina, cabía presumir que esas mutaciones azarosas harían «involucionar» las funciones y capacidades que se buscaba obtener. El resultado más probable sería una avería, el equivalente de una alteración letal: «De modo que, mientras que este sistema resulta extremadamente primitivo, en él se incluye el rasgo de una mutación hereditaria, llegándose al extremo de que una mutación azarosa será probablemente letal, pero podría ser no letal y hereditaria».

Para que las máquinas obtuvieran algo mejor de sus diseños, en esencia una mayor inteligencia, necesitarían que se añadiera un elemento creativo a

sus órganos de estímulos y de fusión. A diferencia de lo que sucede con la evolución biológica, la idea no tendría que esperar millones de años, sino que requeriría que los sistemas madre contuvieran el chispazo prometeico necesario, lo que conduciría de manera más o menos directa a una mejora de los diseños. Aquello era ficción, pensaba Von Neumann. Tal y como comentó ante sus colegas de Princeton, no había ciencia o teorías de ingeniería que pudieran encontrarle un sentido. Von Neumann, que no era ningún ludita, había hecho explotar la «explosión de inteligencia».

Un defecto evidente en las predicciones sobre la explosión de inteligencia que conduzca a la superinteligencia es que ya disponemos de una inteligencia de nivel humano: somos humanos. Siguiendo la lógica de Good, deberíamos ser capaces de diseñar algo por encima de lo humano. Esto no es más que una reafirmación de los objetivos del ámbito de la IA, así que nos hemos metido en un círculo vicioso. Las personas que se dedican a la investigación en IA ya saben que el diseño de artefactos más inteligentes es un misterio, tal y como explicó Von Neumann. Trasladar ese misterio desde nuestra propia inteligencia a la inteligencia imaginaria de una máquina no nos ayuda en nada. Para desentrañar este tema un poco más, pensemos en una investigadora genial de IA a la que llamaremos Alice.

EXPLOSIONES DE INTELIGENCIA, LA IDEA MISMA

Supongamos que Alice es una científica de IA que tiene un vecino tostón llamado Bob. Bob dispone de sentido común, puede leer el periódico y mantener una conversación normal (aunque quizá sea aburrida), así que se encuentra a una distancia sideral por encima de los mejores sistemas de IA surgidos del DeepMind de Google.

Alice trabaja para una empresa emergente (que no tardará en ser adquirida por Google) y quiere construir una IA que sea tan lista como Bob. Ha bosquejado dos sistemas, siguiendo el espíritu de los célebres Sistema 1 y Sistema 2 de Daniel Kahneman.⁴ Se trata de concesiones a la intuición o metáforas que proporcionan un borrador para los tipos de problemas que

habrá que resolver a fin de llegar a la inteligencia artificial general. En el contexto de Alice, los llamaremos Sistema X, de aptitud en tareas bien definidas como la participación en juegos (el ajedrez o el go), y Sistema Y, de inteligencia general. Este último incluye las aptitudes lectoras y conversacionales de Bob, pero también el área más turbia de nuevas ideas y percepciones.

A Bob se le da fatal el ajedrez y, de hecho, su sistema X resulta patético en comparación no solo con el de un sistema como AlphaGo, sino respecto a muchos otros seres humanos. Su memoria a corto plazo es peor que la de la mayoría de la gente; obtiene resultados pobres en las pruebas de inteligencia y le cuesta bastante resolver los crucigramas. En lo que respecta a su sistema Y, su inteligencia general muestra una llamativa falta de interés o de capacidad para el pensamiento novedoso o perspicaz. Bob no es de esos vecinos que reciben demasiadas invitaciones para cenar.

La estrategia de Alice consiste en comenzar diseñando una Máquina de Bob que esté a la altura de su inteligencia. Según su razonamiento, si logra crear una Máquina de Bob, esta podrá diseñar una versión más inteligente de sí misma, lo cual acabará conduciendo a una explosión de inteligencia. Bien, una vez más no olvidemos que diseñar una Máquina de Bob no es tarea sencilla, porque Bob cuenta con un Sistema Y —lo cual significa que ha solucionado el problema del razonamiento y del sentido común, y dispone de capacidades cognitivas generales—. Por ejemplo, podría superar el test de Turing. Y puede leer cuentos infantiles y la sección de deportes y resumirlos. Bob, por tanto, podría machacar a los mejores sistemas de comprensión del lenguaje natural de Google, como la herramienta de búsqueda semántica Talk to Books de Ray Kurzweil. Ese es el motivo por el que Alice está tan emocionada con su proyecto de Máquina de Bob: sería un gran avance para la IA.

La pregunta es: ¿cómo llegar hasta ese punto? El primer enfoque de Alice consiste en maximizar las capacidades del Sistema X de la Máquina de Bob. Le proporciona una memoria de ordenador y acceso a la red a través de Google. Por desgracia, esta versión de la Máquina de Bob no tarda en demostrar la afirmación de Stuart Russell según la cual los superordenadores sin inteligencia real simplemente tardan menos en obtener las respuestas equivocadas.⁵ La Máquina de Bob recuerda las cosas erróneamente y fracasa a la hora de hacer las preguntas adecuadas. Pese a todas las mejoras que va introduciendo en su Sistema X, Alice solo

consigue que la máquina se muestre más competente en recordar y expectorar teorías descabelladas y en realizar declaraciones acerca del mundo con más datos, todos ellos mal utilizados y peor entendidos desde la perspectiva de un Sistema Y. Sí, la Máquina de Bob juega de fábula al ajedrez, pero sus aptitudes para el ajedrez hacen que Alice la considere menos interesante, ya que se da cuenta de que la máquina que ha creado no tiene ninguna posibilidad de diseñar una versión «más inteligente» de sí misma.

En un momento de iluminación, Alice se da cuenta de que Bob no podría diseñar una versión más lista de sí mismo. Por tanto, ¿cómo podría hacerlo la Máquina de Bob? El problema, piensa, es que la optimización del Sistema X no suministra los recursos necesarios al Sistema Y. La Máquina de Bob (igual que el propio Bob) debería cuantificar su inteligencia, evaluar sus limitaciones y la extensión de estas, y a continuación rediseñarse de forma activa a fin de volverse más lista allí donde las cosas cuentan y tienen importancia. ¡Pero sucede que en ese punto es precisamente donde la Máquina de Bob (igual que Bob) se muestra poco inteligente! La Máquina de Bob no puede hacer eso porque su Sistema Y carece de las capacidades de percepción, descubrimiento e innovación. Alice tiene que volver a la casilla de salida.

Entonces, Alice decide que la Máquina de Bob es demasiado estúpida para formar parte del proceso de arranque de una superinteligencia. (En un momento de pánico cerval, se le ocurre que esa lógica pone en peligro la empresa entera de llegar a una superinteligencia, pero se las arregla para reprimir esa ansiedad con rapidez.) Alice decide, en deferencia al fundador de la IA y al entusiasta departamento de mercadotecnia de su compañía, Ultra++, que en su lugar va a concentrarse en el diseño de una máquina tan inteligente como Alan Turing, a la que llamará Máquina de Turing.

Bien, asumiendo que Turing fuera más listo que Alice (aunque ¿quién podría asegurarlo?), ella no puede diseñar directamente esa Máquina de Turing; de hecho ya se había estrellado contra una pared de ladrillos al intentar resolver el diseño de la Máquina de Bob. Decide comenzar con una máquina que sea tan lista como Hugh Alexander —colega de Turing en Bletchley Park y, en una ocasión, campeón de ajedrez en Cambridge—. Hugh Alexander era listo, realmente listo. Jugaba al ajedrez a nivel profesional y, aunque no logró demostrar el mismo nivel de percepción que Turing a la hora de descifrar el código Enigma, sí realizó contribuciones

valiosas a ese empeño y se ganó el respeto de los demás descifradores de Bletchley —que no es poco—. La Máquina de Hugh deberá ser lo bastante lista para averiguar la manera de cablear una Máquina de Turing, ¡y una máquina del nivel de Turing sin duda será lo bastante lista como para volverse más lista aún!

Solo con descargarse en el móvil un código ajedrecístico de StockFish, Alice ya logra mejorar la aptitud de la Máquina de Hugh para el ajedrez (y eso que Hugh fue un campeón en esa disciplina). De manera similar, gracias a una calculadora le otorga a la Máquina de Hugh una capacidad aritmética perfecta, y le añade una memoria de superordenador, así como acceso a toda la información que Google pueda recuperar. El Sistema X es inmejorable y la Máquina de Hugh puede hacer cosas que Hugh Alexander, con toda su inteligencia aparente, nunca pudo llevar a cabo: jugar un ajedrez sobrehumano, realizar sumas imposibles y destacar en muchas otras cosas propias del Sistema X. El problema es que eso también podía hacerlo la Máquina de Bob. De hecho, Alice se da cuenta de que la Máquina de Bob y la Máquina de Hugh son con toda probabilidad equivalentes. Incluso se ve obligada a admitir (tras varios vasos de vino tinto) que abandonar la Máquina de Bob fue innecesario y contraproducente.

Después de otro vaso de vino y un cigarrillo, Alice apaga el móvil para silenciar los molestos mensajes de texto que le están mandando sus colegas de Ultra++ para preguntar por su descubrimiento inminente. La verdad, reflexiona, es que Bob no es el único que no puede diseñar una versión más lista de sí mismo; a Alice le pasa igual. En un momento de lucidez descubre que, cuanto más nos alejamos del Sistema X en dirección al Sistema Y, hacia la percepción y la innovación, más opaco se volverá el diseño. Turing, por ejemplo, podía juzgar su inteligencia en el ajedrez —perdía tanto contra Hugh Alexander como contra Jack Good—, pero no podía evaluar sus propias capacidades de Sistema Y. En un sentido muy real, su intelecto era una caja negra, y bajo ningún concepto podía medir su propia aptitud para el pensamiento original (signifique eso lo que signifique), no solo porque este carece de un plano detallado, por así decirlo, sino porque se asienta en el tiempo, en el plazo de una vida, y aún podría generar ideas nuevas e impredecibles. La inteligencia del Sistema Y de Turing no solo resulta impredecible; en otras palabras, es inexplicable —quizás no ante alguien más listo que Turing (de nuevo: signifique eso lo que signifique), pero sin

duda ante el propio Turing—. Y lo mismo sucede con el lerdo de Bob. ¿Cómo podría Alice, pues, desencadenar una explosión de inteligencia?

Efectivamente, la idea misma de la explosión de inteligencia lleva incorporada una premisa falsa, fácil de exponer a alguien tan ambicioso y perspicaz como Alice, una vez que decide tomársela en serio. Según esa hipótesis, una Máquina de Bob será tan inteligente como Bob. Bien, he aquí una idea: ve a pedirle a Bob que diseñe una versión ligeramente más lista de sí mismo y descubrirás que eso no se encuentra a su alcance. La cualidad esencial de la mente que hace que la IA resulte tan excitante imposibilita a su vez el supuesto lineal de una explosión de inteligencia. «Una vez que lleguemos al nivel de la inteligencia humana, el sistema puede diseñar una versión de sí mismo más inteligente que los seres humanos», se dice con esperanza. Pero es que ya disponemos de una inteligencia de «nivel humano»—somos humanos—. ¿Podemos nosotros hacer algo así? ¿De qué están hablando en realidad los promotores de la explosión de inteligencia?

Es otra manera de decir que los poderes de la mente humana exceden nuestra capacidad para mecanizarla, en el sentido de que se necesita «ampliarla» para pasar del AlphaGo a la Máquina de Bob y a la Máquina de Turing y más allá. La idea misma de la explosión de inteligencia no representa una candidata particularmente buena como Sistema Y para que la IA progrese camino de una inteligencia general.

LOS EVOLUCIONISTAS TECNOLÓGICOS

Muchos entusiastas de la IA se aferran a la tesis de la inevitabilidad (las máquinas superinteligentes van a llegar, hagamos lo que hagamos) porque toca temas evolutivos y, así, de manera muy conveniente, absuelve a los científicos de la responsabilidad individual de realizar hallazgos o desarrollar ideas revolucionarias. La inteligencia artificial se limitará a evolucionar, igual que nosotros. Podemos denominar a esas voces futurólogas y creyentes en la IA dentro de este terreno como «evolucionistas tecnológicos», o ET.

La visión ET es popular entre tecnólogos *new age* como Kevin Kelly, cofundador de *Wired*, quien sostiene en su libro del año 2010 *What Technology Wants* [«Lo que la tecnología quiere»] que la IA no llegará por

obra de un «científico loco», sino simplemente como un proceso evolutivo del planeta, de forma bastante parecida a la evolución natural.⁶ Según su punto de vista, el mundo se está «inteligentizando» (el neologismo es suyo) y cada vez van emergiendo formas de tecnología más y más complejas e inteligentes, y sin un diseño humano explícito.⁷ Es posible que esos pensadores también conciban la World Wide Web, la Red de Extensión Mundial, como un cerebro gigante y en crecimiento. Los seres humanos, según esta perspectiva, devienen un eslabón en la cadena cósmica-histórica que se extiende hacia el futuro y la verdadera IA, momento en que nos quedaremos atrás o seremos asimilados.

La vida orgánica evoluciona con extrema lentitud, pero los ET perciben un progreso tecnológico cada vez más acelerado. Tal y como dice una famosa aseveración de Kurzweil, la tecnología sigue una curva de aceleración y se está volviendo cada vez más complicada, según una ley que él considera discernible a lo largo de la historia: la ley de rendimientos acelerados. Así, la inteligencia de nivel humano y, a continuación, la superinteligencia aparecerán en el planeta en un plazo de tiempo tremendamente corto si lo comparamos con la evolución orgánica. Faltan pocas décadas para que tengamos que enfrentarnos a ellas, quizá solo años.

Esta es una historia de la humanidad pulcra y simple. Nos encontramos en medio de una transición hacia algo diferente, algo que será mejor y más inteligente.

Date cuenta de que esa historia no es comprobable; tan solo tenemos que esperar y ver qué pasa. Si el año que se ha predicho para la llegada de la verdadera IA resulta ser también falso, pues se pronostica otro para dentro de algunas décadas. En ese sentido, la IA es infalsificable y por tanto —según las reglas aceptadas del método científico— carece de rigor, no es una idea científica.

Fijaos en que no estoy diciendo que la IA verdadera sea imposible. Como les gusta señalar a Stuart Russell y a otros investigadores de IA, algunos científicos del siglo xx, como Ernest Rutherford, pensaron que era imposible construir una bomba atómica, pero Leó Szilárd descubrió la manera en que operan las reacciones nucleares en cadena —y lo hizo apenas veinticuatro horas después de que Rutherford diera la idea por muerta—.⁸ Es un buen recordatorio de que no se debe apostar contra la ciencia. Pero piensa que la reacción nuclear en cadena se desarrolló a partir

de unas teorías científicas comprobables. Las teorías acerca de la evolución tecnológica de un poder mental no lo son.

Las declaraciones de Good y Bostrom, presentadas como una inevitabilidad científica, son más bien una concesión a la fantasía: ¡imagínate que esto fuera posible! Y no cabe duda de que sería genial. Y quizá peligroso. Pero imaginar escenarios hipotéticos nos aleja mucho de una discusión seria sobre lo que nos espera.

Para comenzar, una capacidad de superinteligencia general debería estar conectada al resto del mundo de manera que pudiera observar y «hacer conjeturas» de manera más productiva que nosotros. Y, si la inteligencia también es social y situacional, tal y como parece que debe de ser, se requerirá una inmensa cantidad de conocimiento contextual para diseñar algo más inteligente. El problema de Good no es mecánico y restrictivo, sino que más bien atrae hacia su órbita la totalidad de la cultura y la sociedad. ¿Dónde está el plano más simple y remotamente plausible para ello?

En otras palabras, la propuesta de Good se basa, una vez más, en una visión de la inteligencia simplista e inadecuada. Presupone el error original de la inteligencia y le añade otro juego de manos reduccionista: que una inteligencia mecánica individual puede diseñar y construir otra inteligencia mecánica individual superior. Que una máquina pueda situarse en tamaño punto de creación arquimédica parece improbable, por decirlo con suavidad. En realidad, la idea de la superinteligencia es una multiplicación de errores, y representa la esencia del punto al que ha llegado la fantasía en relación con el advenimiento de la IA.

Para ahondar en todo esto tenemos que seguir abriéndonos paso hacia el interior de esa fantasía. Se conoce como la «singularidad», y nos vamos a centrar en ella a continuación.

Capítulo 4

La singularidad, ayer y hoy

En los años cincuenta, el matemático Stanisław Ulam recordaba una vieja conversación que mantuvo con John Von Neumann en la que este comentó la posibilidad de que la humanidad viviera un punto de inflexión tecnológico: «El proceso en aceleración constante de la tecnología ... parece sugerir que nos acercamos a una singularidad esencial en la historia de nuestra raza, tras la cual los asuntos humanos tal y como los conocemos no podrán seguir igual».¹

Lo más probable es que Von Neumann efectuara ese comentario en un momento en que los ordenadores digitales estaban llegando a la escena tecnológica. Pero los ordenadores digitales fueron la última innovación dentro de una larga y en apariencia ininterrumpida secuencia tecnológica.² En los años cuarenta ya había quedado claro que las revoluciones científica e industrial de los trescientos años precedentes habían puesto en movimiento unas fuerzas de un enorme poder simbiótico: los frutos de la nueva ciencia sembraron el desarrollo de la nueva tecnología, que a su vez posibilitó mayores descubrimientos científicos. Por ejemplo, la ciencia nos dio el telescopio, que a su vez permitió grandes avances en astronomía.

Los cambios sociales —rápidos, a veces caóticos y al parecer irreversibles— estuvieron ligados de manera inextricable a los cambios en la ciencia y la tecnología. La población de las ciudades se disparó (con dosis considerables de miseria e injusticia) y, aparentemente de la noche a la mañana, emergieron formas por completo nuevas de organización social

y económica. Las máquinas de vapor revolucionaron el transporte, igual que los motores de combustión interna más adelante. Trenes, tranvías y barcos de vapor abrieron nuevas rutas comerciales y la migración a las ciudades generó nuevas fuerzas de trabajo. Con la invención de la bombilla eléctrica por parte de Thomas Edison, la gente pudo trabajar de noche; en las zonas rurales, quienes padecían insomnio podían ahora leer *El Capital* o *El origen de las especies* cuando el sol ya los hubiera abandonado. La productividad subió como la espuma. La riqueza y la salud crecieron. Lo mismo pasó con la sangre y la violencia. Una secuencia de acontecimientos geopolíticos condujo a la «Gran Guerra», la primera guerra mundial, que introdujo el uso de armas químicas a escala masiva. Y, un par de décadas más tarde, en el mundo de Von Neumann, la amenaza existencial definitiva —la bomba nuclear— se hizo realidad.

La bomba marcó un punto de inflexión histórico, dejó a las claras las posibilidades distópicas inherentes a la innovación tecnológica desenfrenada. Shannon y Turing usaban ordenadores para jugar al ajedrez; científicos como Von Neumann usaron ordenadores para desarrollar armas con las que vaporizar ciudades japonesas. Los ordenadores electrónicos eran grandes y lentos, pero seguían siendo considerablemente más rápidos que los ordenadores humanos para realizar tareas como el cálculo de progresiones numéricas, que Von Neumann y otros usaban para determinar el radio probable de una explosión nuclear según las diferentes cantidades de material fisible.

En este miasma de miedos y posibilidades, Von Neumann fue a plantear la cuestión de una «singularidad». Famoso polímata y científico brillante, Von Neumann disfrutaba de un respeto casi universal entre sus colegas, incluyendo a Alan Turing, y no es de extrañar que sus sugerencias impactaran a Ulam, quien habría de recordarlas décadas más tarde.

No cabe duda de que el matemático Ulam entendió la metáfora de Von Neumann. El de singularidad es un término matemático que indica un punto que se vuelve indefinido —un valor que, pongamos, estalla hacia el infinito—. Von Neumann le preguntó a Ulam si el progreso tecnológico en efecto se acercaría a ese «infinito», sobre el que no se podrían aplicar métodos ni ideas, estrategias ni acciones. Cualquier predicción resultaría imposible. El progreso dejaría de ser una variable conocida (si es que alguna vez lo fue).

En otras palabras, Von Neumann sugirió a Ulam una escatología, un posible final de los tiempos. Un par de décadas más tarde, Good creyó

haber encontrado el mecanismo para ello: el ordenador digital.

Vernor Vinge, científico informático de la UCLA y ganador del premio Hugo, introdujo en 1986 la «singularidad» en la informática y, más concretamente, en la inteligencia artificial con su libro de ciencia ficción *Naufragio en el tiempo real*.³ En un artículo técnico posterior para la NASA, Vinge canalizó las ideas de Good:

En un plazo de treinta años dispondremos de los medios tecnológicos para crear una inteligencia superhumana. Poco después, la era de la humanidad habrá acabado ... Creo que es justo definir este acontecimiento como una singularidad. Se trata de un punto en el que debemos descartar nuestros modelos y una nueva realidad gobierna. A medida que vayamos acercándonos a ese punto, su sombra se cernirá con un tamaño cada vez mayor sobre los asuntos humanos, hasta que la noción se convierta en un lugar común. Sin embargo, cuando suceda al fin, quizá siga representando una gran sorpresa y un enigma aún mayor.⁴

Vinge, el científico informático, contó con la compañía de varios profesionales. A finales de los ochenta, Raymond Kurzweil, científico informático del MIT, futurólogo y emprendedor, se había convertido en el *bulldog* de la IA al expandir la idea de la singularidad en la ciencia de la cultura popular con una serie de publicaciones, comenzando con *The Age of Intelligent Machines* [«La era de las máquinas inteligentes»], de 1990, y su secuela de 1998, *La era de las máquinas espirituales*. En su *best seller* de 2005 se mostró aún más confiado: *La Singularidad está cerca*.

Kurzweil afirmaba que la innovación tecnológica, presentada en un gráfico histórico, es exponencial. La innovación se acelera, desde una visión histórica, como una función de tiempo. En otras palabras, el tiempo que transcurre entre una innovación tecnológica de primer orden y la siguiente no hace más que menguar. Por ejemplo, el papel apareció en el siglo II, y la imprenta tardó otros mil doscientos años en llegar —la imprenta de Gutenberg apareció en Europa en 1440—. Pero la computación apareció en los años cuarenta (en los años treinta, si contamos su tratamiento matemático), e internet —una innovación bastante capital— apareció menos de treinta años después. ¿Y la IA? Según la lógica de Kurzweil, la inteligencia de nivel humano en un ordenador se encuentra a pocas décadas de distancia, quizá menos—. Las curvas de crecimiento exponencial nos sorprenden.

Kurzweil etiquetó esa idea como «ley de rendimientos acelerados» (LOAR en sus siglas inglesas) y la usó como premisa en un debate cuya conclusión fue que la IA completamente humana llegaría en 2029 y a

continuación, a través de una serie de procesos de carga automática hacia máquinas más inteligentes, la superinteligencia lo haría hacia el año 2045.⁵

La superinteligencia marcaba el punto sin retorno, el lugar donde la senda del progreso desaparece en lo desconocido, en la singularidad. Se trata del punto de intersección, donde las máquinas, y no la gente, se convierten en los seres más inteligentes del planeta.

Como es sabido, Kurzweil considera que ese proceso es completamente «científico»; cita para ello la LOAR (aunque la LOAR no sea ninguna ley) y, en gran medida, se apoya en su propio entusiasmo y en sus credenciales como experto informático e inventor (Kurzweil ayudó a desarrollar la tecnología de conversión de texto a voz, lo que condujo a sistemas modernos como Siri).

Turing. Good. Vinge. Las ideas acerca de un cambio radical posibilitado por los avances informáticos ya estaban en el aire. A todas luces Kurzweil les proporcionó un mapa de carreteras. Como sucede con tantos otros sujetos obsesionados con el tema de la IA, su prosa es más papista que el papa:

Estamos entrando en una nueva era. Yo la llamo «la Singularidad». Se trata de una fusión de la inteligencia humana y la inteligencia de la máquina que va a crear algo más importante que la suma de sus partes. Se trata del proceso evolutivo más vanguardista de nuestro planeta. Se podría argumentar con firmeza que, de hecho, se trata de la vanguardia de la evolución de la inteligencia en general, porque no existen indicios de que haya sucedido en ningún otro lugar. Para mí, es la esencia misma de la civilización humana. Es algo que forma parte de nuestro destino y del destino de la evolución de cara a seguir progresando aún a mayor velocidad, y para que el poder de la inteligencia crezca de manera exponencial. Pensar en ponerle fin —pensar que los seres humanos ya estamos bien así— forma parte de un recuerdo bonito pero inapropiado de lo que fue la raza humana. La raza humana es una especie que ha experimentado una evolución cultural y tecnológica, y la naturaleza de esa evolución consiste en acelerarse, y en que su potencia crezca de manera exponencial, y de eso es de lo que estamos hablando. Su estadio siguiente consistirá en amplificar nuestra propia capacidad intelectual con los resultados de nuestra tecnología.⁶

No obstante, para ser sinceros hemos de decir que, en el momento en que Kurzweil se subió al tren con tanto entusiasmo —de hecho, décadas antes de 1980—, el trabajo en pos de una IA científica ya había extinguido la esperanza de una marcha inexorable hacia la superinteligencia. La investigación en IA y su desarrollo se habían revelado, en una palabra, difíciles.

En los años setenta, Hubert Dreyfus, filósofo del MIT y mosca cojonera de la IA, había publicado ya una influyente refutación de esa disciplina en

cuanto ejemplo clamoroso de lo que el filósofo húngaro Imre Lakatos llamó «programa de investigación degenerativo».⁷ Dreyfus se puso estupefacto, pero algo de razón llevaba, como los propios científicos informáticos saben demasiado bien. La disciplina sufría un revés tras otro, con aquellos intentos bien financiados y aquellas declaraciones grandilocuentes sobre máquinas inteligentes que se iban quedando cortos de manera constante (y a menudo dramática). Los laboratorios de investigación del MIT, Stanford y demás lugares se encontraron con una sucesión en apariencia interminable de disyuntivas, dificultades, confusiones y fracasos absolutos. Por ejemplo, en los años cincuenta se pensó que, dedicándole el esfuerzo de investigación y los dólares suficientes, se podría resolver el problema de la Máquina Traductora Completamente Automatizada de Alta Calidad. En los sesenta, tras una sucesión de fracasos, la inversión gubernamental en traducción había desaparecido. La esperanza de construir robots dotados de sentido común (pongamos que de la capacidad para entender el inglés y hablarlo) también se había evaporado —o al menos se había visto drásticamente reducida por aquella oleada de decepciones tempranas—. Los sistemas conversacionales que debían pasar el test de Turing de manera realista lograban burlar a los interrogadores humanos solo gracias a sus trucos y sus engaños —no a una comprensión real de la lengua—, problema que continúa haciendo mella hoy en día en los trabajos en pos de un lenguaje natural en la IA. Esta aparecía inevitablemente en las notas de prensa y en las charlas sobre el futuro, pero no sucedía lo mismo cuando la atención se centraba en el trabajo real de los laboratorios de investigación. Programar una máquina que fuera inteligente de verdad resultó ser difícil. Muy difícil.

Mientras la idea de la singularidad iba ganando fuerza en la cultura popular, los científicos de la IA seguían adentrándose en problemas de ingeniería en apariencia interminables. Y el cielo no acababa de caer. La singularidad no se acercaba. La ficción popular de Vinge seguía siendo eso: ficción.

Un examen más minucioso de la IA revela una brecha vergonzosa entre el progreso real que han realizado los científicos informáticos que se dedican a ella y las visiones futuristas que estos y otras personas disfrutaban describiendo. En 1950, Turing propuso la comprobación de una pregunta: ¿pueden las máquinas ser tan listas como las personas? Good, Vinge, Kurzweil y compañía han contestado a ella con un sí rotundo, sin tomarse

en serio la verdadera naturaleza de los problemas con los que se encuentra el trabajo en ese terreno.

Esa brecha resulta muy ilustrativa.

En concreto, el fracaso de la IA a la hora de seguir un proceso sustancial sobre los aspectos más complicados de la comprensión de un lenguaje natural sugiere que las diferencias entre mente y máquina son más sutiles y complicadas de lo que Turing se imaginó. Y, si la historia de la IA nos ha de servir de guía, estas representan una dificultad profunda para la disciplina.

Nos centraremos en ellas a continuación.

Capítulo 5

La comprensión del lenguaje natural

La inteligencia artificial como disciplina oficial comenzó, con buenos auspicios, en 1956, en la ahora famosa Conferencia de Dartmouth. Entre los asistentes hubo celebridades como Shannon, de Bell Labs (teoría de la información), Marvin Minsky, de Harvard (matemáticas), Herbert Simon (destacado economista de Carnegie Mellon), John McCarthy, George Miller (psicólogo de Harvard conocido por su trabajo sobre la memoria humana) y John Nash (el matemático laureado con el Nobel al que retrató la popular película de 2001 *Una mente maravillosa*).

Durante la conferencia, McCarthy, que por entonces estaba en Dartmouth, pero muy pronto iba a aceptar un puesto en Stanford para desarrollar el nuevo ámbito de la ciencia informática, acuñó el término «inteligencia artificial» y dio así nombre oficial al proyecto moderno de diseñar vida inteligente. En 1816, la joven y precoz Mary Shelley había comenzado a trabajar en su obra maestra, *Frankenstein*. Ciento cuarenta años más tarde, los científicos reunidos en Dartmouth consideraron el montaje de un nuevo «moderno Prometeo», que no tardaría en irrumpir ante la opinión pública.

La disciplina cayó en la hipérbole desde el minuto uno. Las actas de la conferencia lo decían todo:

Proponemos que, a lo largo de dos meses del verano de 1956, diez hombres lleven a cabo un estudio sobre la inteligencia artificial en el Dartmouth College de Hanover, Nuevo Hampshire. Ese estudio debe tomar como base la conjetura de que todos los aspectos del aprendizaje o de cualquier otro rasgo de la inteligencia se pueden describir por principio con la precisión

necesaria como para que una máquina pueda imitarlos. Se intentará averiguar cómo lograr que las máquinas utilicen el lenguaje, establezcan abstracciones y conceptos, resuelvan tipos de problemas que ahora están reservados a los seres humanos y se perfeccionen a sí mismas. Creemos que se puede realizar un avance notable en uno o más de estos problemas si un grupo de científicos cuidadosamente escogido trabaja de manera conjunta en ellos durante un verano.¹

La agenda de Dartmouth fue simple: investigar la naturaleza de las capacidades cognitivas (el pensamiento), diseñar programas que las reprodujeran e implementar y comprobar su desempeño en las nuevas computadoras electrónicas. Tal y como los participantes de Dartmouth dejaron claro, en el verano de 1956, gracias a diez investigadores armados con el conocimiento de sus respectivos campos científicos, esperaban un «avance significativo» de cara a diseñar una inteligencia humana en una máquina.

Mientras trabajaban en RAND, Herbert Simon y Allan Newell diseñaron a finales de los años cincuenta unos programas de IA que parecían cumplir las promesas optimistas de la Conferencia de Dartmouth. El programa de IA Logic Theorist, y más tarde el Solucionador General de Problemas, se sirvieron de una simple búsqueda heurística para demostrar teoremas de lógica tradicional y resolver puzzles de base lógica con claros pasos computacionales. Los programas funcionaron y la IA pareció destinada a desentrañar con rapidez los secretos de la inteligencia humana, tal y como habían declarado los organizadores de Dartmouth.

Los éxitos iniciales de Simon y Newell no tardaron en animar a otros investigadores para que se marcaran objetivos más ambiciosos. Turing ya había señalado el objetivo final del programa una década antes con su versión del juego de la imitación: el test de Turing. Los científicos de Dartmouth también pensaron que programar una máquina para que entendiera el inglés o cualquier otra lengua natural representaría una declaración de victoria para la IA. Los investigadores llevaban mucho tiempo pensando que la comprensión de un lenguaje natural sería «IA completo», jerga que habían tomado prestada de las matemáticas para indicar que en el momento en que las computadoras dominaran un lenguaje natural habrían alcanzado la inteligencia general, y serían por consiguiente capaces de pensar y actuar como los seres humanos. En los años sesenta, pues, el objetivo de la IA era la traducción automática —el trasvase completo y automatizado de textos de una lengua, como podría ser el ruso, a otra, como el inglés—. La IA «iba a por todas».

LA COMPRENSIÓN DEL LENGUAJE NATURAL

Cuando la IA se centró en la comprensión del lenguaje natural, sus fieles irradiaron una seguridad en el éxito inminente de la empresa, digno de la tradición iniciada en Dartmouth. Herbert Simon, quien acabaría ganando el prestigioso premio A. M. Turing y más tarde el premio de Economía Conmemorativo de Alfred Nobel, anunció en 1957 que «ahora mismo hay en el mundo máquinas que piensan, que aprenden y que crean». En 1965 predijo que, hacia 1985, «las máquinas serán capaces de realizar cualquier trabajo que haga un hombre». Marvin Minsky también declaró en 1967 que «en el lapso de una generación, el problema de crear una “inteligencia artificial” se habrá solucionado de manera sustancial».²

Pero los investigadores no tardarían en descubrir que el de la traducción automática era un partido completamente diferente. Partieron del supuesto simplista de que se podía entender el lenguaje a partir del análisis de las palabras de textos amplios (llamados «corpus») utilizando técnicas estadísticas, pero pronto se demostró que estaban equivocados. Las computadoras posibilitaron la traducción automática, pero los resultados no fueron de buena calidad. Ni siquiera los programas que operaban en dominios específicos, como el de la literatura biomédica, eran inmunes al error, y esos errores eran a menudo estúpidos y vergonzosos.

Los investigadores de la traducción automática respondieron expandiendo su estrategia y explorando métodos para «disecionar» las frases, para hallar su estructura sintáctica, usando las nuevas y potentes gramáticas «transformacionales» que había desarrollado un joven lingüista del MIT que pronto iba a obtener fama mundial: Noam Chomsky. Pero la labor misma de extraer análisis correctos de los textos de lenguaje natural se reveló muchísimo más difícil y compleja de lo que nadie se hubiera imaginado. Aparecieron problemas que, en retrospectiva, deberían haber resultado evidentes. Entre ellos se encontraba la ambigüedad en el sentido de las palabras (cuando una palabra como «banco» presenta diferentes significados posibles), la dependencia contextual deslocalizada (cuando el significado de la palabra depende de otras palabras del texto o discurso que no se encuentran próximas a ella) y otros asuntos relacionados con la referencia (anáfora), la metáfora y la semántica (el significado). Tal y como lo definió el filósofo y científico cognitivo Jerry Fodor, la IA se había

metido en una partida de ajedrez tridimensional pensando que iba a jugar al tres en raya.³ A mediados de los años sesenta, el National Resource Council de Estados Unidos invertía millones de dólares en los trabajos sobre traducción automática de numerosas universidades del país, pero, por decirlo suavemente, los sistemas de ingeniería que comprendieran o que incluso fingieran comprender los textos de lenguaje natural con éxito eran bastante escasos.

El investigador del MIT Yehoshua Bar-Hillel, en su momento partidario feroz y entusiasta de la traducción completamente automatizada, fue el primero en dar la voz de alarma. Hizo algo más que eso, de hecho, con una serie de informes oficiales para el NRC, acompañados de unas ahora famosas notas a pie de página, en los que explicó con todo lujo de detalle los profundos problemas a los que se enfrentaba la disciplina.⁴ Sus informes tuvieron un efecto sísmico sobre la comunidad investigadora. Había precisado con exactitud el obstáculo contra el que se estrellaba la traducción automática, y este era irritantemente «filosófico»: la escasez de un supuesto sentido común o «conocimiento del mundo» —el conocimiento del mundo real—. Pensemos en esta frase simple del inglés: «*The box is in the pen*», que la máquina podría traducir al castellano por «La caja está en el bolígrafo». Bar-Hillel explicó que esta frase servía para confundir a todos los sistemas automatizados, sin importar su grado de sofisticación, porque estos carecían de un conocimiento simple y real del mundo. El conocimiento sobre los tamaños relativos de bolígrafos y cajas capacita a los seres humanos para eliminar la ambigüedad de ese tipo de frases de manera casi instantánea. Reconocemos con rapidez que lo más probable es que el segundo elemento de la frase no sea un instrumento de escritura, sino un recinto para niños pequeños o animales, puesto que «pen» significa tanto «bolígrafo» como «redil» o «parque infantil». Y todo se vuelve más evidente con un poco de contexto adicional, como en el ejemplo de Bar-Hillel: «El pequeño John buscaba la caja de sus juguetes. Al final la encontró. La caja estaba en el parque infantil. John se puso muy contento». Pero los sistemas automatizados, que carecen de ese conocimiento, se enfrentan a una tarea misteriosa y en apariencia imposible.

Tal y como señaló Bar-Hillel, a las computadoras no se les podía suministrar ese conocimiento del mundo —al menos no de manera directa, ingenieril— porque «el número de datos que nosotros, los seres humanos, conocemos es, en un sentido creativo y fértil, infinito».⁵ Sin querer había

descubierto que las personas sabemos mucho más de lo que cualquiera podría haberse imaginado —lo cual representaba lo opuesto a una solución rápida y simple para la IA—. Y eran los hechos cotidianos, de sentido común, en apariencia mundanos, acerca de la vida diaria, los que confundían a los más sofisticados sistemas automatizados. Cualquier hecho en apariencia normal podía volverse relevante en el curso de una traducción y hacer que los volúmenes de conocimiento necesario y abierto a los que podían acceder las computadoras en tiempo real o casi real, y que les infundían la capacidad cognitiva para seleccionar hechos relevantes contra ese contexto abierto (posiblemente infinito), parecieran inútiles. Tal y como concluyó Bar-Hillel en su célebre informe de 1966 para el NRC, la idea de que las computadoras pudieran ser programadas con el conocimiento del mundo de los seres humanos era «una completa quimera y difícilmente merece un debate mayor».⁶

En otras palabras, la traducción automática se había atascado con unos resultados que estaban a años luz de ser traducciones completamente automáticas de calidad elevada (y siguen siendo así, aunque la calidad haya mejorado). Por consiguiente, se mantenía el patrón. La IA se había excedido con la propaganda y, a resultas del fracaso de la investigación en traducción para estar a la altura de las promesas realizadas, el NRC retiró su financiación tras haber invertido más de veinte millones de dólares en investigación y desarrollo, una suma enorme en aquel momento. Como resultado de aquella debacle, los investigadores de IA perdieron sus trabajos, se destruyeron carreras y la IA como disciplina tuvo que volver a la casilla de salida.⁷

Durante los años setenta y ochenta, las tentativas por controlar o resolver el «problema del conocimiento del sentido común» dominaron los esfuerzos de la investigación en IA. No obstante, a principios de la década de 1990, la IA seguía sin contar con una estrategia o respuesta novedosa para su problema medular científico —y filosófico—. Japón había invertido millones en la Quinta Generación, un proyecto de perfil alto que aspiraba a alcanzar éxitos en el terreno de la robótica, y también fracasó de manera harto espectacular. A mediados de aquella década, la IA volvió a sus cuarteles de invierno —sin confianza en las promesas de sus investigadores, sin resultados para demostrar a sus detractores que estaban equivocados y sin financiación—. Entonces llegó la red.

LA RED DE EXTENSIÓN MUNDIAL

La aparición de la World Wide Web, la Red de Extensión Mundial, estimuló el resurgimiento de la IA por una simple razón: los datos. De repente, la disponibilidad de conjuntos masivos de datos, y sobre todo de corpus textuales (páginas web), debidos al esfuerzo conjunto de millones de nuevos usuarios de la red, insufló vida a las viejas y «superficiales» estrategias de estadística y reconocimiento de modelos. De repente, lo que solía ser superficial se volvió adecuado y comenzó a funcionar. Los algoritmos de aprendizaje supervisado, como las redes neuronales artificiales (redes neuronales, para abreviar), los árboles de decisión y los clasificadores bayesianos existían desde hacía décadas en los laboratorios universitarios. Pero, sin conjuntos amplios de datos, aún no habían revelado su potencial en problemas de interés como el reconocimiento facial, o la clasificación de textos, o el correo basura, o la detección de fraudes. De repente, esos métodos parecían estar llenos de promesas sin fin —entre ellas, la de crear aplicaciones rentables en el mundo real que iban a provocar una nueva ola de atención y de financiación hacia la IA.

Así nacieron los macrodatos, el *big data* (aunque el término llegó un poco más tarde). Con el cambio de siglo, los enfoques tipificados por algoritmos de aprendizaje como las redes neuronales y los modelos gráficos —supuestamente superficiales, ascendentes, empíricos o impulsados por datos— pasaron a amparar vastas oportunidades tanto en el terreno de la investigación en IA como en el de sus aplicaciones comerciales. Se desarrollaron nuevos métodos —incorporando modelos de Márkov ocultos, modelos de entropía máxima, campos aleatorios condicionales y clasificadores de margen amplio, como máquinas de vectores de soporte— que pasaron a dominar con rapidez la investigación pura y aplicada en la IA. En apariencia, de la noche a la mañana surgió una ciencia completa de análisis numéricos y estadísticos, basada en la optimización de los métodos de aprendizaje que trabajaban con macrodatos. Las universidades pusieron en marcha proyectos de comprensión y procesamiento del lenguaje natural. Encontraron la manera, por ejemplo, de extraer nombres y otros modelos de las páginas web (una habilidad conocida como reconocimiento de entidades), de desambiguar palabras polisémicas (con múltiples sentidos) como «banco», de realizar tareas específicas de la red como la ordenación y

recuperación de páginas web (siendo su ejemplo más famoso el PageRank de Google, que Larry Page y Sergey Brin desarrollaron en la década de 1990 como estudiantes de posgrado en Stanford), de clasificar noticias y otras páginas web por temática, de filtrar el correo basura en el correo electrónico y de ofrecer recomendaciones de productos espontáneas en sitios comerciales como Amazon. La lista no tiene fin.

El trasvase entre los enfoques lingüísticos y basados en reglas y los métodos impulsados por datos o «empíricos» pareció liberar la IA de aquellos primeros días de trabajo tormentoso en la traducción automática, cuando los problemas en apariencia interminables sobre la captura de significados y el contexto azotaron las labores de ingeniería. De hecho, la traducción automática se acabó resolviendo gracias a un grupo de investigadores de IBM que usaron un enfoque estadístico (es decir, de base no gramatical) que, en esencia, no dejaba de ser una aplicación ingeniosa de los trabajos iniciales de Claude Shannon sobre la teoría de la información. Se le conoce como el enfoque del «canal ruidoso» y considera las frases de un lenguaje de origen (pongamos que el francés) y un lenguaje de destino (pongamos que el inglés) como un intercambio de información en el que las malas traducciones constituyen una forma de ruido —lo que conduce a que la tarea del sistema consista en reducir el ruido del canal de traducción entre las frases de origen y de destino—. La idea funcionó, y las máquinas comenzaron a generar traducciones útiles sirviéndose del enfoque estadístico liderado por los laboratorios de investigación de IBM, muchísimo más simple pero con un uso intensivo de datos.

ÉXITO... O NO

El éxito de sistemas contemporáneos como Google Translate ante el problema antaño desconcertante de la traducción automática se publicita a menudo como una prueba de que la IA acabará triunfando si se le proporcionan el tiempo suficiente y las ideas adecuadas. La realidad invita a poner los pies en el suelo.

Por un lado, resulta que ciertos problemas sobre la comprensión de un lenguaje natural se pueden abordar con enfoques estadísticos o de aprendizaje automático; pero, por otro, las preocupaciones originales de

Bar-Hillel y otros acerca de la semántica (el significado) y la pragmática (el contexto) se han revelado fundamentadas. La traducción automática, que parecía representar un problema relacionado con el lenguaje natural de muy difícil solución, pudo realizarse de manera adecuada gracias a análisis estadísticos simples que contaran con corpus amplios (conjuntos de datos) en lenguajes diferentes. (Y asumamos que la traducción automática aún no presenta una calidad demasiado elevada; vendría a ser solo «razonablemente buena».) Esto no demuestra que la inteligencia de las máquinas haya crecido de manera impactante a la hora de entender el lenguaje natural, sino que la traducción automática es un problema mucho más sencillo de lo que se pensó en un primer momento.

Una vez más, siguen existiendo grandes problemas para que los ordenadores comprendan el lenguaje. Una manera simple de verlo consiste en recuperar el test de Turing y replanteárselo a la luz de la historia de la IA, con numerosos y en su mayoría infructuosos esfuerzos para resolver los problemas que la acompañan, o incluso para realizar algún avance sustancial. Lo más probable es que futurólogos como Nick Bostrom, junto con la mayor parte de la comunidad científica de la IA, deseen que el público se olvide del test.

No se trata —como se dice a veces— de que el test sea fallido o sirva de poco. Simplemente es demasiado difícil.

EL TEST DE TURING

Si se las observaba a diez mil metros de distancia, en efecto pudo parecer que las computadoras iban ganando inteligencia mientras la IA progresaba desde su génesis con los primeros trabajos de Turing y el inicio de la conferencia de Dartmouth. Sin duda, las computadoras pasaron a tener procesadores más potentes y una mayor capacidad de almacenaje, y eso condujo de manera inevitable al desarrollo de unos problemas bien definidos. A la hora de jugar al ajedrez, las computadoras no sobrepasaban la capacidad de las personas más cualificadas, pero tampoco dejaban de mejorar. Programas como DENDRAL, que analizaba estructuras químicas, y MYCIN, que a veces ofrecía diagnósticos médicos bastante buenos, dejaron a las claras que los métodos de la IA eran relevantes para diversos

problemas que por lo general requerirían de una inteligencia humana elevada. Tal y como ya hemos visto, en un primer momento la traducción automática fue un fracaso, pero comenzó a ofrecer rendimiento con diferentes enfoques posibilitados por la disponibilidad de grandes conjuntos de datos (fue la antesala de numerosos éxitos con los macrodatos y la ciencia de datos en los 2000). Los sistemas de IA, cada vez más potentes y sofisticados, comenzaron a realizar una labor de zapa con todo tipo de tareas de procesamiento de lenguaje natural, como la de generar análisis sintácticos de oraciones de lenguaje natural y etiquetar partes del discurso o entidades (personas, organizaciones, lugares y cosas así).⁸

Sin embargo, seguía sin alcanzarse el objetivo original de Turing para la IA: que lograra pasar su famoso test. Para la IA, comprender una conversación no era una tarea bien definida; se trataba de toda una aptitud. Aunque existió alguna esperanza de éxito con sistemas tempranos como el SHRDLU de Terry Winograd, que apilaba y manipulaba bloques con un brazo robótico a través del análisis sintáctico de órdenes como «Coloca el bloque A encima del B, y coloca ambos encima del C», los investigadores no llegaban a nada con los sistemas conversacionales abiertos. Los avances se basaban en engaños, en estratagemas ingeniosas de los diseñadores de sistemas para burlar a los interlocutores humanos fingiendo ignorancia o adoptando respuestas evasivas y tácticas diversas cuando las preguntas sacaban al programa de su elemento. El programa ELIZA de Joseph Weizenbaum, quizá el sistema trucado de mayor fama, podía generar la ilusión de un diálogo humano imitando los comentarios de un psicoterapeuta rogeriano a alguno de sus pacientes. He aquí una muestra, con las respuestas de ELIZA en cursiva:

- Todos los hombres son iguales.
- En qué sentido*
- Siempre nos están chinchando por un motivo u otro.
- Se te ocurre un ejemplo específico de lo que dices*
- Bueno, mi novio me ha hecho venir aquí.
- Tu novio te ha hecho venir aquí*
- Dice que estoy deprimida la mayor parte del tiempo.
- Lamento oír que estás deprimida*
- Es cierto. No soy feliz.
- Creo que venir aquí te ayudará a ser más feliz*

—Necesito ayuda, al menos eso parece ser cierto.⁹

Pero esos programas, según los parámetros de la verdadera IA, eran fraudes evidentes.

Como también lo fue el programa que en múltiples ocasiones burló momentáneamente a los jueces del premio Loebner, una puesta en escena del test de Turing con carácter anual. La prensa se apresuró a declarar que al fin se había superado el test, y siguió una breve pero intensa celebración a bombo y platillo del avance de la IA ante su más famoso desafío. El bot conversacional —llamado Eugene Goostman— se presentó como un muchacho ucraniano de trece años, identidad diseñada para que los jueces anticiparan (y excusaran) sus errores gramaticales y sus lagunas de conocimientos generales. El bot también contestó a algunas preguntas no con respuestas directas, sino con distracciones y evasivas en forma de comentarios aleatorios y chascarrillos digresivos. ¿Engañó el bot por momentos a algunos de los jueces? Sí. ¿Superó en realidad el test de Turing de manera razonable? No.¹⁰

Aunque durante estos últimos años se ha podido dominar (en cierto modo) la traducción automática gracias a los amplios volúmenes de textos traducidos a diferentes idiomas que hay en la web, el test de Turing sigue representando una frustración perpetua para la IA. El fantasma de Bar-Hillel continúa acosándonos.

Capítulo 6

De la IA como tecnología *kitsch*

En 1980, el escritor de origen checo Milan Kundera escribió su obra maestra, *La insoportable levedad del ser*. Se trata de una historia de amor con el telón de fondo de la invasión de Checoslovaquia en 1968 por parte de la Unión Soviética. Kundera escribió acerca de los escritores y artistas que se suicidaron tras ser acosados con calumnias de forma incansable por la policía secreta soviética, cuyos elementos se habían insertado en el tejido social, intelectual y cultural de Praga. Muertos y desacreditados, los intelectuales praguenses sufrían a continuación una vergüenza más (aunque póstuma): el elogio repugnante, durante sus funerales, por parte de miembros y oficiales del partido soviético, que daban fe de la devoción por el Estado que el finado había sentido a lo largo de toda su vida. La propaganda soviética los conducía a la muerte y, acto seguido, esa misma propaganda presentaba sus existencias como si las hubieran sacrificado con nobleza para promover unas ideas contra las que, de hecho, se habían manifestado tanto en público como en privado. Se decía que habían amado aquello que más odiaban.

La propaganda soviética era implacable, pero no era ni furiosa ni estúpida. Tenía un propósito concreto, y ese propósito consistía en purgar el país de aquellas expresiones profundas y elevadas (y opositoras) acerca del sentido del país, de su gente y de su vida. Los soviéticos purgaron Praga, y Checoslovaquia entera, de su historia compartida, de sus tradiciones y de su sentido de lo que resultaba precioso, aquello por lo que

valdría la pena luchar. Tras silenciar a los librepensadores, como quien pinta una pared después de pulirla, los soviéticos serían libres para imponer su visión del mundo sin enfrentarse a una oposición seria u organizada. El relato de Kundera es un recuento incisivo y a menudo trágico sobre el valor de la vida humana y la manera en que una creencia o ideología en particular puede intentar —pero nunca conseguir por completo— ofuscar y encubrir todo aquello que resulta importante para el individuo y la sociedad. Kundera definió la cultura que los soviéticos endosaron al pueblo checo derrotado como *kitsch*.

LA TECNOLOGÍA *KITSCH*

«*Kitsch*» es una palabra alemana que, si bien hoy en día se aplica a obras de arte o motivos decorativos cursis o chabacanos, en su origen implicaba un sentimentalismo y melodrama exagerados en cualquier ámbito. Los errores de inteligencia en el núcleo de la cosmovisión de la IA —sus creencias; es decir, no la ciencia— han dado pie a una forma moderna y particularmente perniciosa de *kitsch*. El sueño de un superordenador inteligente no es como la propaganda soviética, y nadie nos obliga a creer en el advenimiento de las máquinas. Pero ambos sí comparten la idea básica de reemplazar debates difíciles y complejos sobre el individuo y la sociedad con historias tecnológicas que, al igual que la cultura soviética, reescriben ideas viejas con abstracciones peligrosamente unidimensionales.

El significado y el uso de la palabra *kitsch* han cambiado con el paso del tiempo. La definición original en alemán difiere de algún modo del significado que pretendo explorar aquí, pero hay dos ingredientes esenciales de ese sentido inicial que deberían ofrecer la suficiente claridad a mi afirmación. Primero, *kitsch* implica la simplificación de unas ideas complicadas. Tiene que haber una historia simple que se pueda contar. Segundo, ofrece soluciones sencillas que barren con la emoción las preguntas y la confusión que la gente experimenta hacia sus problemas vitales, en vez de abordar esas preguntas con un debate serio y exploratorio. Así, encontramos un ejemplo perfecto de *kitsch* en la ensoñación de que algún día un androide impresionante y dotado de superinteligencia remodelará la sociedad humana y sus viejas tradiciones e ideas, y que

entraremos en una nueva era, por fortuna libre de las antiguas discusiones acerca de Dios, la mente, la libertad, la buena vida y demás. Unas máquinas hermosas (o unas máquinas dotadas de una bella inteligencia), como la Ava de la película de ciencia ficción de 2015 *Ex Machina*, a la que dio vida Alicia Vikander, acabarán con la dura realidad de la existencia humana. Eso es *kitsch* al estilo tecnológico. Como la propaganda soviética, puede horrorizarnos o apaciguarnos, pero nos proporciona una historia nueva que se escribe sobre la anterior y que vuelve innecesario lo que era verdad antes, de modo que la antigua realidad desaparece.

Pese a todas sus contribuciones a la ciencia y la ingeniería, Alan Turing hizo posible la génesis y el crecimiento viral del *kitsch* en la tecnología al comenzar equiparando la inteligencia con la resolución de errores. Más adelante, Jack Good agravó el error de la inteligencia de Turing con su muy discutida idea de ultrainteligencia y la propuesta de que la aparición de las máquinas inteligentes implicaría de manera necesaria la llegada de las máquinas superinteligentes. Una vez que la imaginación popular hubo aceptado la idea de las máquinas superinteligentes, la reescritura del sentido de la existencia humana, su significado e historia se podría efectuar dentro de los parámetros de la informática y de la tecnología.

Pero las máquinas ultrainteligentes son una fantasía, y pretender lo contrario da aliento a ese proceso indeseable que conduce a la tecnología *kitsch*, por lo general de dos maneras igual de superficiales. En un extremo escuchamos un relato de IA apocalíptico o aterrador, como una de esas historias que se cuentan alrededor del fuego en un campamento. En el otro extremo nos encontramos con la IA utópica o soñadora, que es igual de superficial e inmerecida. Si nos tomamos en serio una u otra forma de la IA *kitsch*, acabaremos en un mundo definido tan solo por la tecnología.

Es un tema que retomaré más tarde, porque expone el problema fundamental de la IA futurista. Tal y como dice Nathan, el genial científico informático de *Ex Machina*: «Llegará el día en que las inteligencias artificiales vuelvan la vista atrás y nos vean tal y como nosotros vemos los esqueletos fósiles de las llanuras africanas. Un simio erecto que vivió rodeado de polvo, con una lengua en bruto y herramientas sin pulir, listo para extinguirse». La verdad es que no está claro que vaya a haber algún ordenador capaz de volver la vista atrás. Ese sentimiento popular requiere de una inmersión profunda en el sentido de la existencia, de la vida, de la consciencia y de la inteligencia, y en las diferencias que existen entre

nosotros y la informática y sus numerosas tecnologías. Lo *kitsch* evita que tratemos de resolver la naturaleza humana y demás empeños filosóficos de gravedad. Y no debería ser así, como bien sabía Kundera.

Lo *kitsch* hunde sus raíces, por lo general, en un sistema de pensamiento más amplio. Para los comunistas, este fue el marxismo. En el mito de la inevitabilidad, es la tecnociencia. Nuestra visión tecnocientífica del mundo es una herencia directa de la obra de Auguste Comte.

NUESTRA CONDICIÓN TECNOCIENTÍFICA

Probablemente, el primer pensador que explicó y desarrolló de manera completa la cosmovisión de la tecnociencia fue Auguste Comte.¹ Comte, filósofo del siglo XIX reconocido popularmente por haber sentado los fundamentos de la sociología como un campo de estudio científico, desarrolló y expuso la teoría del positivismo. Según esa visión, los únicos fenómenos que existen son aquellos observables y científicos; la religión y la filosofía serían imaginarios. Comte explicitó, primero, su idea de que la mente humana se dirige hacia la verdad, igual que la sociedad en su conjunto, a través de unas etapas que se inician con los pensamientos religioso y filosófico, y que progresan hacia el pensamiento científico. Y, segundo, explicó que la tecnociencia acabaría por crear un cielo en la tierra al posibilitar la comprensión de la naturaleza de todas las cosas (ciencia), y al servirse a continuación de ese conocimiento para desarrollar tecnologías que alarguen de manera notable nuestras vidas, que las mejoren y que las vuelvan más valiosas.

La explicación de Comte sobre el poder transformativo de la tecnociencia acabó dando pie a la convicción de que la religión y la Iglesia en particular podrían ser reemplazadas por una «religión de la humanidad» plenamente secular, que no creería en ningún dios y que estaría anclada con firmeza en las ciencias y en la realidad material. En el momento en que Comte escribió todo eso, durante el siglo XIX, había suficientes pruebas sobre el poder del pensamiento humano tanto para descubrir leyes científicas como para crear tecnologías útiles y potentes, de modo que la tecnociencia arraigó en el centro mismo de la mente moderna.

No obstante, la teoría de Comte provocó recelos desde un primer momento. Nietzsche, por ejemplo, lamentó que la idea de persona se viera restringida y limitada de esa manera. La tecnociencia podía ayudarnos a vivir más tiempo, pero no podía hacernos más sabios. La idea del héroe, o de una persona con dones y virtudes extraordinarios y merecidos, no entraba en la visión de Comte, que en esencia había reemplazado el debate tradicional acerca de la persona con el debate sobre el progreso de la ciencia y, en especial, de la tecnología.²

El materialismo de Comte también sugirió a otros pensadores una disminución, en vez de la expansión, de las posibilidades humanas. Desde el este, en Rusia, el escritor Dostoyevski protestó desdeñoso contra el «flagelo» creciente de la creencia absoluta en el materialismo y el cientificismo —la visión según la cual el científico es el único conocimiento real—, y lo hizo con una prosa que reflejaba el escepticismo de otros pensadores e incluso su miedo ante la velocidad con la que iba ganando predominancia el pensamiento tecnocientífico. Tal y como dijo en *Memorias del subsuelo*: «Nuestro propio deseo, voluntario y libre; nuestro propio capricho, aun el más alocado; la fantasía desatada hasta rayar en lo extravagante...; he ahí en qué consiste la ventaja pasada por alto, el interés más principal, que en ninguna clasificación se incluye y que manda a paseo todos los sistemas y teorías».^{*a 3}

Dostoyevski, Nietzsche y otros apuntaban al ideal de la persona plena, pero Comte hablaba de un ideal externo a nosotros: la tecnociencia y su progreso. El problema consistía en que, como bien sabía Comte, la visión de un futuro tecnocientífico representaba también una declaración profunda y cargada de sentido sobre la naturaleza de la persona. En efecto, Comte afirmó que las concepciones tradicionales de la persona —como algo único porque lo creó Dios, o como alguien en busca de la sabiduría (no solo del conocimiento tecnológico), tal y como sostenían los filósofos griegos— habían pasado a ser irrelevantes en virtud de los éxitos científico y tecnológico. Su filosofía tecnocientífica fue una apostilla a la esencia y las posibilidades de la naturaleza humana. Era una postura radical, y los pensadores iconoclastas, que no se dejaron engañar por el monstruo de la tecnociencia, acertaron al desafiar las ideas propuestas por Comte (y compañía).⁴

EL TRIUNFO DEL *HOMO FABER*

La tecnociencia triunfó durante el siglo xx, pero también siguió generando respuestas escépticas. Hannah Arendt, la filósofa que se hizo famosa por su frase sobre «la banalidad del mal», en referencia a los juicios nazis de Núremberg, afirmó que la tecnociencia de Comte —que a mediados del siglo xx ciertamente no había perdido fuerza como idea filosófica— equivalía nada más y nada menos que a una redefinición de la naturaleza humana misma.⁵ Arendt apuntó hacia la comprensión clásica de los seres humanos como *Homo sapiens* —literalmente, «hombre sabio»— y al foco histórico sobre la sabiduría y el conocimiento en vez de las aptitudes técnicas, y argumentó que abrazar la tecnociencia como cosmovisión implicaba redefinirnos como *Homo faber* —«hombre constructor».

En términos griegos, el *Homo faber* es una persona que cree que la *téchne* —el conocimiento de un oficio o acerca de la construcción de algo, la raíz de la tecnología— define aquello que somos. Esa visión *faberiana* de la naturaleza humana se ajusta a la perfección no solo a la idea decimonónica de Comte sobre una tecnociencia utópica, sino a la obsesión del siglo xx por construir tecnologías cada vez más potentes y que culminó con el proyecto grandioso de, en efecto, construirnos a nosotros mismos: la inteligencia artificial. Ese proyecto no tendría sentido si las nociones tradicionales sobre el significado de la humanidad hubieran permanecido intactas.

Arendt sostuvo que el cambio sísmico desde la sabiduría y el conocimiento, y hacia la tecnología y la construcción, implicaba una comprensión de nosotros mismos limitante y peligrosa en potencia, que garantizaría no solo que el desarrollo tecnológico prosiguiera desenfrenado, sino que viéramos cada vez más los éxitos tecnológicos como declaraciones valiosas acerca de nosotros mismos. En otras palabras, estábamos reduciendo nuestra propia valía a fin de aumentar, más allá de lo sabio o razonable, nuestra estimación hacia las maravillas que se podían construir con las herramientas de la tecnociencia.

Los comentarios en un principio crípticos de Von Neumann sobre esa aceleración del avance tecnológico que nos acerca a la «singularidad» ganan claridad a la luz de la posición de su contemporánea Arendt. Aunque Von Neumann, científico y matemático, no ahondó (hasta donde sabemos)

en esas afirmaciones, estas reflejan a la perfección la insistencia de Arendt en la significación profunda de la tecnociencia para nosotros y para nuestro futuro —para lo que los filósofos de la tecnología llaman «la condición humana»—. A Comte quizá le pareciera perverso que la tecnología pudiera acelerarse al extremo de escapar a nuestro control, pero en ninguno de sus textos se puede encontrar el menor indicio de lo que iba a apuntar Arendt (y otros); que, al abogar por la tecnociencia como respuesta humana a los problemas humanos, nos estamos involucrando en el proyecto de redefinir nuestra comprensión de nosotros mismos. El giro hacia la *téchne* en vez de, pongamos, la *episteme* (conocimiento de los fenómenos naturales) o la *sapientiae* (la sabiduría relacionada con los valores humanos y sociales) dificulta la forja de una idea valiosa de la unicidad humana. (Al fin y al cabo, incluso las abejas son constructoras; de colmenas, en su caso.)

Situar la *téchne* en el centro del debate también posibilita la visión de la persona como algo que se puede construir, ya que implica que su máxima consecución como tal es una capacidad superior para fabricar tecnologías cada vez más avanzadas. Una vez emprendida esa ruta, el trayecto hacia la inteligencia artificial es corto. Y ahí aparece la conexión evidente con los errores de la inteligencia que cometió primero Turing y que se hicieron extensivos a Jack Good y compañía hasta llegar al día de hoy: el triunfo definitivo del *Homo faber* como especie consiste en construirse a sí mismo. Y ese es precisamente, por supuesto, el objetivo declarado de la IA. Explorar la posibilidad de que el proyecto sea un éxito o no requerirá que nos sumerjamos en las aguas profundas de la comprensión de nuestra propia naturaleza.

PARA COMPLETAR EL PUZLE

La tecnociencia se inició con la revolución científica, y pocos siglos después la mayor parte de la teoría científica moderna ya había ocupado su lugar. Con raras excepciones —siendo una de ellas, evidentemente, el desarrollo de la teoría cuántica y de la relatividad durante el siglo xx—, el conocimiento científico avanzó a medida que se ponían en marcha la mayoría de las teorías de la física. El conocimiento científico era como un puzle de piezas teóricas que conformaban una imagen del mundo y del

universo. La física de Newton, la electrodinámica de Maxwell, las teorías del trabajo y la termodinámica de Carnot y otros..., todas esas formas de conocimiento científico encajan entre sí y ofrecen una imagen unificada del mundo. La teoría de la evolución de Darwin, en la década de 1850, y los sucesivos descubrimientos geográficos y arqueológicos añadieron nuevas hipótesis y detalles. (Por supuesto, esas teorías eran debatidas y comprobadas, y algunas se revelaban erróneas o se revisaban.) El alcance de las posibilidades de la teoría científica, pues, iba menguando de manera extraña —como cuando uno trabaja con un rompecabezas, y con cada pieza que encaja en su lugar el número de elecciones que quedan se va volviendo cada vez más limitado.

En contraste, la innovación tecnológica estalló. Tal y como ha señalado Ray Kurzweil, la innovación tecnológica se acelera. Una invención no limita aquello que puede venir a continuación, sino que posibilita que haya cada vez más invenciones. En otras palabras, la tecnología parece evolucionar. No la ordenamos igual que la teoría. En cambio, hacemos crecer los desarrollos tecnológicos apilándolos unos encima de los otros de manera en apariencia interminable. La aceleración de la evolución tecnológica implica simplemente que el plazo de tiempo entre una innovación tecnológica capital y la siguiente no deja de menguar, históricamente, así que la distancia entre, pongamos, la invención de la imprenta y la llegada de la computadora es muy amplia en comparación con, pongamos de nuevo, la distancia entre la computadora e internet. La fusión entre ciencia y tecnología resulta, pues, complicada, y la palabra misma «tecnociencia» sugiere que, a medida que las cosas vayan progresando, la ciencia se acomodará y la tecnología continuará evolucionando —y lo hará, tal y como dice Kurzweil, de manera exponencial.

Así que el término mismo «tecnociencia» demuestra la complejidad y el carácter impredecible de nuestro mundo. No todas las áreas de la actividad humana siguen los mismos parámetros de crecimiento; ninguna de ellas puede tenderse a lo largo de otra, como si se tratara de una plantilla. Que la inteligencia humana se parezca a la inteligencia de la máquina —o no— es algo que aún está por verse. La cuestión de la IA debería ser una invitación no a ignorar los problemas filosóficos, sino a pelearse con ellos. Y la tecnociencia, entendida como una afirmación sobre nosotros mismos, es en definitiva una simplificación terrible, que representa (entre otras cosas) la

introducción de lo *kitsch* en la corriente de los asuntos complejos y difíciles de esta vida.

Capítulo 7

Simplificaciones y misterios

Poco antes de que Turing publicara en 1950 su «Maquinaria computacional e inteligencia», el psicólogo conductista B. F. Skinner publicó la novela de ciencia ficción *Walden dos*.¹ En ella, Skinner hace que sus personajes afirmen que el libre albedrío es una ilusión, y que el comportamiento de las personas se puede controlar de manera externa, desde su entorno. Si alguien (pongamos que un científico) altera el entorno, el comportamiento de esa persona en ese entorno se alterará también.

En un sentido trivial, eso es verdad. Si un déspota le niega comida, seguridad y oportunidades de trabajo a la gente, esa gente se volverá infeliz. Podemos predecir ese tipo de cambios. Lo que Skinner quería decir, en cualquier caso, es que la persona se encuentra completamente determinada por impulsos externos (por estímulos, en sus propias palabras).

De hecho, la noción skinneriana de la persona como una «caja negra» es la misma idea básica que Turing tenía en la cabeza. En las cajas negras, tratamos la salida del sistema como una función de su entrada —el funcionamiento de su interior no se llega a describir—. Skinner afirmó en *Walden dos* que se podría construir un mundo perfecto —una utopía— tratando a la gente como si fueran cajas negras, proporcionándoles ciertos impulsos físicos (estímulos) para obtener ciertos resultados (respuestas). Mientras tanto, Turing especulaba con la idea de que el ser humano era equivalente operativamente a una máquina compleja, y para demostrarlo

sugirió construir una máquina, proporcionarle impulsos y examinar sus respuestas.

Por desgracia, esa manera de pensar prescindía de muchas cuestiones de importancia, y hoy parece evidente que hemos heredado sus errores. Mientras que la teoría del condicionamiento operante de Skinner —o conductismo, como se le llamó después— iba a ser motivo de gran controversia más tarde, durante el siglo xx, la interdisciplinaria «revolución cognitiva» que la reemplazó pasó a tratar la inteligencia como un conjunto de meros cálculos internos. Esa idea, apuntalada por una filosofía denominada «teoría computacional de la mente», que afirma que la mente humana es un sistema de procesamiento de información, no deja de respaldar la confianza teórica en el triunfo futuro de la IA.

Aquí lo mejor es ser claros: equiparar la mente humana con un ordenador no es una actitud científica, sino filosófica.

EL DISPARATE DE LA PREDICCIÓN

Tal y como señala Stuart Russell, en la búsqueda de la inteligencia artificial no deberíamos apostar en contra del «ingenio humano».² Pero, de manera similar, no deberíamos realizar predicciones optimistas (o catastrofistas) sin una sólida base científica.

A los expertos e incluso (o en especial) a los científicos les encanta hacer predicciones, pero la mayoría de ellas son erróneas. En su excelente libro *Future Babble*, Dan Gardner documenta la tasa de acierto de las predicciones realizadas en los ámbitos de la historia y la geopolítica, hasta llegar a las ciencias.³ Y descubre que los teóricos —expertos con grandes visiones de futuro basadas en la teoría particular que ellos mismos respaldan— tienden a realizar predicciones peores que la gente pragmática, que ve el mundo como algo complejo, que no acaba de encajar con ninguna teoría en concreto.

Gardner se refiere a la clase de los expertos y a los pensadores pragmáticos como «erizos» y «zorros» respectivamente (toma prestados los términos del psicólogo Philip Tetlock, quien a su vez le había copiado la terminología a Isaiah Berlin). De la misma manera que el erizo escarba su madriguera, los expertos erizos escarban en una idea. De forma inevitable

llegan a pensar que esa idea plasma la esencia de todo, y esa fe alimenta el proselitismo consiguiente. Marx fue un erizo infatigable.

Los zorros ven la complejidad y el carácter incalculable de los asuntos del mundo, y o bien evitan realizar predicciones audaces o predicen con mayor seguridad (y quizá inteligencia) que las cosas no cambiarán de la manera que pensamos. Para el zorro, el negocio de la predicción es casi una temeridad, porque la verdad es que no podemos saber lo que surgirá de la compleja dinámica que existe entre la geopolítica, la política interna de un estado (pongamos, ¿quién ganará unas elecciones?), la ciencia y la tecnología. Tal y como nos advirtió Lev Tolstói, novelista del siglo XIX, las guerras tienen lugar por motivos que no podemos encajar en los planes de batalla.

Algunos científicos de IA son famosos por haberse mostrado vulpinos en sus predicciones sobre esa disciplina. Mira a Yoshua Bengio, profesor de ciencias de la informática en la universidad de Montreal, Canadá, y uno de los pioneros del aprendizaje profundo: «No soy yo quien pueda decirlo — responde a la pregunta de si podemos esperar una IA de nivel humano—. No tiene sentido ni sirve de nada adivinar una fecha porque no tenemos ni idea. Lo único que puedo decir es que no sucederá en un plazo de pocos años».⁴

Ray Kurzweil ofrece una respuesta más propia de un erizo: la IA de nivel humano llegará en 2029. Invoca la «ley» de rendimientos acelerados para hacer que sus predicciones parezcan científicas y no deja de encontrar pruebas de que ha acertado con todos los supuestos avances que han tenido lugar hasta la fecha.⁵

A veces, los filósofos tienen la virtud de pensar con claridad sobre un problema porque no les estorba ningún celo concreto que pueda ligarse a los profesionales de ese terreno (los que aún deseen filosofar). Por ejemplo, Alasdair MacIntyre, en su ya clásico *Tras la virtud*, señaló las cuatro fuentes fundamentales de impredecibilidad que hay en el mundo. En concreto, su exposición sobre la «innovación conceptual radical» tiene una relevancia directa en las preguntas acerca de la llegada de una IA de nivel humano. MacIntyre recuerda el argumento contra la posibilidad de predecir inventos que realizó Karl Popper, filósofo de la ciencia del siglo XX:

En algún momento del Paleolítico, tú y yo estamos hablando del futuro y yo predigo que antes de que transcurran diez años alguien habrá inventado la rueda.

—¿La rueda? —preguntas—. ¿Y eso qué es?

Procedo a describirte la rueda, buscando las palabras, sin duda con dificultad, que por primera vez sirvan para decir lo que serán el aro, los radios, el buje y quizá el eje. Acto seguido hago una pausa, traspuesto.

—Pero es que nadie podrá inventar la rueda, porque la acabo de inventar yo.

En otras palabras, la invención de la rueda no se puede predecir. Y es que una parte necesaria a la hora de predecir un invento consiste en decir lo que la rueda es, y decir lo que la rueda es implica inventarla. Es sencillo apreciar la manera en que este ejemplo puede generalizarse.

Cualquier invento, cualquier descubrimiento cuya esencia consista en elaborar un concepto radicalmente nuevo, no se puede predecir, ya que una parte necesaria de esa predicción es la elaboración actual del concepto mismo cuyo descubrimiento o invención solo debía tener lugar en el futuro. La idea de predecir una innovación conceptual radical es en sí misma una

incoherencia conceptual.⁶

En otras palabras, sugerir que nos encontramos en la «senda» de una inteligencia artificial general cuya llegada se pueda predecir presupone que no hay ninguna innovación conceptual de camino, una postura con la que ni siquiera se mostrarían de acuerdo los científicos de IA más convencidos del advenimiento de la inteligencia artificial general y que estén dispuestos a realizar predicciones, como Ray Kurzweil. Como mínimo, todos sabemos que para que un supuesto sistema de inteligencia artificial general alcance una capacidad todavía desconocida para comprender el lenguaje natural, antes deberá darse el invento o el descubrimiento de un componente racional y generalizador. Esto cuenta sin duda a modo de ejemplo de «innovación conceptual radical», porque aún no tenemos ni idea de lo que será, ni de su aspecto siquiera.

La idea de que podemos predecir la llegada de la IA esconde por lo general una premisa, en grado diverso de aceptación, según la cual el éxito en los sistemas de IA restrictivos, como el de los juegos, facilitará una ampliación que lleve a la inteligencia general, de modo que se puede trazar con cierta confianza la línea predictiva que conducirá de la inteligencia artificial a la inteligencia artificial general. Se trata de una mala suposición, tanto para alentar los avances en el terreno hacia la inteligencia artificial general como para la lógica de los argumentos en favor de una posible predicción.

Las predicciones acerca de descubrimientos científicos quizá se puedan entender mejor como una forma de complacencia mitológica; en efecto, la certeza acerca de la llegada de la inteligencia artificial general solo puede morar en el reino de lo mítico, sin las trabas que le suponen las dudas de Popper o de MacIntyre o de cualquier otra persona.

No toda la mitología sobre la IA es mala. Esta mantiene con vida el anhelo arquetípico por crear vida e inteligencia, y puede abrir ventanas hacia la comprensión de nosotros mismos. Pero, cuando se disfraza de ciencia y de certeza, el mito confunde al público y frustra a los investigadores fuera del ámbito mitológico, conscientes de que aún hay que resolver grandes obstáculos teóricos. «No tenemos ni idea», en palabras de Bengio. Aunque cuente con el apoyo de todas las pruebas y sea cierto, esto resulta extremada y deprimentemente pesimista para los mitólogos.

Sin embargo, los obstáculos no son siempre infranqueables, e, incluso cuando lo son —cuando nos vemos obligados a reconocer ciertos límites—, a continuación quedamos liberados para encontrar una manera diferente de alcanzar nuestro objetivo, o, según el impulso, para formular unos objetivos completamente nuevos. La historia de la ciencia está llena a reventar de ejemplos en los que el descubrimiento de un punto muerto condujo a un nuevo avance. Werner Heisenberg descubrió el principio de incertidumbre mientras intentaba solucionar las consecuencias de la nueva física cuántica. El principio sostiene que es imposible aislar la posición y la cantidad de movimiento de una partícula subatómica de manera simultánea. Eso plantea unos límites fundamentales sobre nuestra capacidad para predecir los movimientos individuales de las partículas en el campo subatómico (porque «ver» la posición de la partícula requiere hacerla chocar con un fotón, lo que también tiene el efecto de desviarla). Aunque no sea más que una limitación, el principio de incertidumbre se ha revelado tan valioso como fructífero a la hora de comprender la mecánica cuántica. Por ejemplo, no habríamos podido aspirar a construir un ordenador cuántico si antes no hubiéramos comprendido la naturaleza de la incertidumbre.

Hay muchos más ejemplos. El móvil perpetuo fue una obsesión de los siglos XVIII y XIX; su órbita atrajo a las mejores y más brillantes mentes. Los avances en las teorías del trabajo y de la termodinámica jubilaron ese sueño..., pero, de paso, ampararon un progreso inmenso en la comprensión de la energía y el movimiento. Aceptar la complejidad —y las complicaciones— nos conduce más lejos que las simplificaciones excesivas.

UN EXTRAÑO (PERO OPORTUNO) ARGUMENTO POR PARTE DE MICHAEL POLANYI

Una posibilidad que ha surgido en el debate sobre la IA es que alcancemos la inteligencia general pero no podamos anotarla —es decir, programarla— porque en muchas áreas importantes será una caja negra para nosotros. Esto nos lleva a hablar de Michael Polanyi.

Polanyi, químico y filósofo influyente en su día pero ahora poco conocido, afirmó a mediados del siglo xx que los símbolos con los que la anotamos solo capturan una parte de la inteligencia —los usos del lenguaje que él llamó «articulaciones»—. Polanyi anticipó muchos de los dolores de cabeza que los sistemas de IA han causado a sus diseñadores; de hecho, en sus trabajos tardíos rechazó explícitamente que las máquinas pudieran capturar la totalidad de la inteligencia humana por razones derivadas del carácter incompleto de las articulaciones.

Polanyi sostuvo que las articulaciones obvian por necesidad componentes «tácitos» de la inteligencia —aspectos del pensamiento que no se pueden describir con precisión anotando símbolos—. ⁷ (La red neuronal que construimos también es un sistema de símbolos.) Por ejemplo, eso sirve para explicar que ciertas habilidades o artesanías, como la gastronómica, no se puedan dominar simplemente tras leer unas recetas. Hacemos cosas, pero eso no significa que podamos programar todo lo que hacemos: piensa en escribir un programa que escriba una novela del orden del *Ulises* de James Joyce. Ese programa carecería de sentido. En su lugar, escribiríamos la novela directamente (si fuéramos Joyce).

Los escritos de Polanyi llegaron en un mal momento para sugerir una visión opuesta a la IA, ya que la disciplina había despegado en los años cincuenta a bombo y platillo. Su defensa del conocimiento tácito fue recogida más adelante por Hubert Dreyfus con su ya comentado ataque contra la IA; quizá a causa de su tono, en ocasiones demasiado tendencioso, los comentarios de Dreyfus se convirtieron en un pararrayos de réplicas y, al menos en un primer momento, no se ganaron a los pensadores más populares de la IA. (Por desgracia, también declaró que un sistema de IA nunca podría vencer a un gran campeón de ajedrez.) ⁸

Pero la posibilidad de que no todo lo que sabemos se pueda anotar representa un problema recurrente para la IA, ya que implica que los programadores están buscando la cuadratura del círculo. Estos escriben programas específicos (o programas de análisis de datos, lo cual sigue siendo específico) a los que se les escapa algo relacionado con nuestra mente. Las ideas de Polanyi sugieren que la mente y la máquina presentan diferencias fundamentales y que equiparar la una con la otra conduce a una simplificación de nuestras ideas sobre la mente. Si la mente —o al menos la inteligencia general— debe ser tratada como algo codificable o anotable, tenemos que simplificar la «mente» misma para que buena parte del debate contemporáneo cobre sentido.

EL REGRESO DE LOS ZORROS

A principios de la década de 2000, la IA se hallaba integrada solo por zorros. La disciplina estaba experimentando uno de sus inviernos perennes y la mayoría de los mitólogos se habían escondido. Ray Kurzweil continuaba promocionando su visión cargada de confianza y los teóricos de la IA clásica, como Doug Lenat, seguían dedicados a sus teorías favoritas, buscando la piedra Rosetta de la IA. Pero los altibajos en apariencia interminables habían desgastado buena parte de la disciplina, hasta el extremo de que muchos se sentían incómodos ante la idea de usar siquiera la etiqueta de IA en nuestras investigaciones. La IA se había convertido en una mala publicidad. (Quizá parezca extraño hoy en día, pero es cierto.) El debate no tardó en apuntar hacia el esoterismo de los algoritmos específicos, como las «máquinas de vectores de soporte» y la «máxima entropía», ambos enfoques de aprendizaje automático. Los científicos de la IA clásica los habían descartado como «superficiales» o «empíricos», porque los enfoques estadísticos que usaban datos no utilizaban el conocimiento y no se les daban demasiado bien ni el razonamiento ni la planificación (si es que no se les daban directamente mal). Pero, cuando la red comenzó a proporcionar esos datos tan necesarios, los enfoques comenzaron a revelarse prometedores.

La «revolución» del aprendizaje profundo se inició hacia 2006, con los primeros trabajos de Geoff Hinton, Yann LeCun y Yoshua Bengio. En 2010,

Google, Microsoft y otras compañías *big tech* ya usaban redes neuronales en aplicaciones de consumo masivo como la del reconocimiento de voz, y en 2012 los *smartphones* de Android contenían la tecnología de las redes neuronales. A partir de ese momento y hasta 2020 (cuando escribo esto), el aprendizaje profundo ha sido el martillo de todos los clavos en lo que a los problemas de la IA se refiere: problemas que cabría abordar «desde cero», como los juegos y el reconocimiento de voz y los datos de imagen, ocupan la mayor parte de la investigación y de los dólares comerciales en IA.

Con el despegue del aprendizaje profundo, la IA (y la conversación en torno a la IA) tomó vuelo también. Los erizos regresaron y, de manera predecible, los medios atizaron las llamas de aquel nuevo futurismo. Pero algo extraño ha estado pasando últimamente en la IA. Reparé en ello en 2018, a través de voces más escépticas, y en 2019 resultaba ya inconfundible. Están volviendo los zorros.

Muchos mitólogos (con algunas excepciones, pocas, pero notables) no son expertos, como Elon Musk o el fallecido astrofísico Stephen Hawking o incluso Bill Gates. Aun así, ayudaron a generar buena parte del bombo de los medios sobre la IA —principalmente, el bombo del aprendizaje profundo—, que alcanzó su cenit hace poco (alrededor de 2015, año arriba, año abajo). Ahora, no obstante, vuelve a ser cada vez más habitual que se hable de limitaciones —es lo que hace, por ejemplo, Gary Marcus, científico cognitivo y fundador de la compañía de robótica Robust.AI, coautor en 2019, junto con el científico informático Ernest Davis, de *Rebooting AI: Building Artificial Intelligence We Can Trust*⁹ [«El reinicio de la IA: La construcción de una inteligencia artificial en la que podamos confiar»]. Marcus y Davis presentan con argumentos de peso la idea de que la disciplina vuelve a estar sobredimensionada, y de que el aprendizaje profundo tiene sus límites; que se necesitará algún avance fundamental para alcanzar una IA de inteligencia general. En 2017, el científico de IA Hector Levesque (colega de Davis, sobre quien hablaremos con mayor detenimiento más adelante) firmó una provechosa polémica sobre la IA moderna que tituló *Common Sense, the Turing Test, and the Quest for the Real AI*¹⁰ [«Sentido común, el test de Turing y la búsqueda de la IA real»]. Cuando publiqué el artículo «Questioning the Hype about Artificial Intelligence» [«El furor de la inteligencia artificial, bajo la lupa»] en *The Atlantic*, en 2015, las reacciones fueron en general desdeñosas.¹¹ En la actualidad hay más voces críticas, y entre ellas se encuentran las de

numerosos líderes reconocidos de la IA, que en efecto cuestionan ese bombo y platillo.

Sigue siendo raro oír argumentos reflexivos sobre la imposibilidad de una IA real, y eso se debe al mismo motivo por el que la gente se abstiene de realizar predicciones al respecto: porque se desconoce el futuro de la IA. Pero, en lo cultural y en lo psicológico, la disciplina parece haber entrado en una fase de relajación, en la que se advierte a los principiantes y al público expectante de que habrá que recorrer un largo camino para llegar a la inteligencia general. Esta tendencia tiene una importancia capital, porque el mito es un faro emocional mediante el cual navegamos por el tema. Es expansionista, acoge a todas las visitas: conceptos como la consciencia, emociones como la agresividad o el amor, instintos como el sexo y otros ingredientes propios de la mente y de los seres vivos. Pero el nuevo debate «científico» es, más o menos, una narrativa sobre las posibles extensiones de la IA débil hacia una generalización cada vez mayor, de cuyo alcance quedan fuera ideas a gran escala como la de la consciencia. Se pasa de lista, quizás... el mito es lo que motiva el interés de todo el mundo. De otro modo, la cosa se queda en formas de tecnología cada vez más potentes por todas partes, una tendencia que, como enseguida percibimos, tiene dos caras.

HACIA UNA SUPERINTELIGENCIA SIMPLIFICADA

Los errores de inteligencia que ayudaron a forjar nuestro mundo computacional simplificado también han regresado bajo una apariencia moderna. Stuart Russell, coautor de la guía introductoria a la IA definitiva junto con Peter Norvig, de Google, asegura en *Human Compatible: Artificial Intelligence and the Problem of Control* [«Compatible con la humanidad: La inteligencia artificial y el problema del control»], de 2019, que la inteligencia no implica nada más que alcanzar objetivos —nos ofrece una definición que incluye no solo a los seres humanos, sino a los delfines y las hormigas, a la bacteria *E. coli* y a los ordenadores—. ¹² Además, quiere que se jubile el test de Turing porque ha pasado a ser irrelevante. (Al

parecer, mantener una conversación normal y corriente no es un objetivo digno.)

El test de Turing no resulta útil para la IA —escribe— porque se trata de una definición informal y extremadamente accidental, que depende de la complicadísima y muy desconocida peculiaridad de la mente humana, resultado tanto de la biología como de la cultura. No existe la posibilidad de «vaciar» la definición y comenzar desde cero a crear máquinas que superen el test de manera comprobable. En su lugar, la IA se ha concentrado en el comportamiento racional, [y por tanto] la máquina será inteligente siempre y cuando sus acciones tengan como resultado probable la consecución de lo que quiere, según lo que haya percibido.¹³

Cuesta discutirle a Russell su definición de inteligencia, que lo abarca todo: desde el momento en que Einstein «alcanzó» su «objetivo» al reimaginar la física como algo relativo, hasta la margarita que vuelve su rostro hacia el sol. Pero el rechazo de Russell hacia el test de Turing parece restrictivo y legalista en exceso, ya que el espíritu del test consiste tan solo en que las máquinas que puedan comprender y utilizar un lenguaje natural deben ser inteligentes. En términos prácticos, no deberíamos esperar gran cosa de Siri y demás asistentes personales activados por la voz si nunca se van a enterar de lo que nos están diciendo, así que desdeñar el test no parece muy inteligente. (Si una Siri de próxima generación alcanza algún día el punto de mantener conversaciones ilimitadas y comunes y corrientes con su dueño humano, el test de Turing regresará en cuanto gran «sueño de la IA» realizado al fin. Pero ay...)

Russell, un reconocido experto en IA y profesor de ciencias de la informática en la universidad de California, Berkeley, también se quita de encima el problema de la consciencia: «En el terreno de la consciencia, la verdad es que lo ignoramos todo, así que no pienso decir nada». A continuación nos asegura que «dentro de la IA no hay nadie trabajando en crear máquinas conscientes, ni sabría nadie por dónde comenzar, y ningún comportamiento tiene la consciencia como prerequisite». Pero, de todos modos, sí dice algo —bastante— acerca de la consciencia:

Supón que te doy un programa y te pregunto: «¿Representa algún tipo de riesgo para la humanidad?». Tú analizas su código y, en efecto, al ejecutarlo, este traza y lleva a cabo un plan cuyo resultado será la destrucción de la raza humana, tal y como un programa de ajedrez traza y lleva a cabo un plan cuyo resultado será la derrota de cualquier ser humano que se enfrente a él. Ahora supón que te digo que ese código, al ejecutarlo, también genera una forma de consciencia automática. ¿Haría cambiar eso tu predicción? En absoluto. No representa la menor diferencia. Tu predicción acerca de su comportamiento sigue siendo la misma, porque se basa en el código. La verdad es que todos esos argumentos hollywoodienses sobre máquinas que cobran

consciencia de manera misteriosa y comienzan a odiar a los seres humanos no han entendido nada: lo que importa es la aptitud, no la consciencia.¹⁴

Pero quizá sea Russell el que no ha entendido nada, porque la mitología de las máquinas que «cobran vida» conforma en realidad el alma de los sueños acerca de una futura IA. Si informáramos a la gente que llega a la sala de cine para ver un pase de *Ex Machina* de que el verdadero sueño de la IA consiste en crear superordenadores mecánicos, «sin luces en su interior», que nos ayuden (a nosotros y a nuestros enemigos) a alcanzar objetivos, quizá se quedarían muy poco impresionados. Russell parece sugerir que un sistema algorítmico, animado de la manera adecuada por unos módulos aún desconocidos para obtener una inteligencia general, representará el éxito definitivo de la IA. Los verdaderos erizos entienden aquello que a los zorros cautelosos se les escapa: que la IA incluye la ciencia y el mito, y que su fascinación perdurable en la mentalidad popular implica que se trata de una piedra de toque psicológica y cultural.

Ray Kurzweil ha sostenido en todo momento que, sea lo que sea la consciencia, las máquinas dispondrán de ella a patadas —y que será más rica y «mejor» que la nuestra—. En 1999, su peán al mito llevó el apropiado título de *La era de las máquinas espirituales* (y lo de «espirituales» lo dijo en serio, como si los ordenadores superinteligentes fueran a tener experiencias conscientes e incorpóreas). Kurzweil insiste sabiamente en que el test de Turing representa una cota de referencia adecuada para la IA: «Para superar el test has de ser inteligente». Incluso le desagrade el concepto de IAG, que cada vez se usa más para hablar de manera específica sobre la inteligencia artificial general, ya que, tal y como afirma (correctamente), «el objetivo de la IA siempre ha consistido en alcanzar grados de inteligencia cada vez mayores, y en última instancia llegar a los niveles humanos de inteligencia».¹⁵

El del deseo sexual podría ser un tema apropiado para la IA, en cuanto prueba con tornasol de la inteligencia —sobre todo si se trata de un elemento básico de nuestros esfuerzos y anhelos, y para alcanzar objetivos diversos, como parece ser en la vida—. *Ex Machina* resulta prácticamente shakespeariana en su combinación de tensión sexual, consciencia, explotación y liberación —y todo ello en una especie de test de Turing—. El novelista reconvertido en realizador Alex Garland aborda la singularidad mientras nos ofrece la historia de una androide superinteligente que planea su huida de la situación de esclavitud a la que la mantiene sometida un

científico loco (otro tema cargado de contenido): su solitario inventor, Nathan (interpretado por Oscar Isaac). Presumimos que el objetivo de Ava es superar un test de Turing interactuando con un invitado de Nathan, Caleb Smith (Domhnall Gleeson), y resultando completamente convincente como «humana», pese a que él ha sido informado de que Ava es un androide y lo ve, además—. Se trata del test definitivo, dice Nathan.

Pero Ava tiene sus propias ideas (a la mierda nuestros objetivos) y planea su huida al mundo salvaje, más allá de los confines del centro de investigación de Nathan. Cuando Ava huye al fin, dos humanos están (o estarán) muertos. Al salir ve los colores, unos colores gloriosos, una demostración para el espectador de que está «viva» y consciente de veras.¹⁶ La hemos visto comprender y utilizar el inglés de manera tan efectiva que ha reducido a los dos hombres a un estado de confusión desesperada y derrota. Aquí tenemos, pues, una representación a voz en grito del mito, que muestra la idea futurista y fundamental de un próximo cruce de caminos en el que las máquinas adelantarán a la humanidad, y punto. Ava es más lista, más astuta, más espiritual, y está más viva que sus equivalentes humanos.

La visión de Garland es pura mitología —también es una gran historia humana, que captura temas arquetípicos (la liberación, el bien y el mal, y la sexualidad) a través de la lente de una tecnología futura—. Se da la ironía de que el éxito de *Ex Machina* se debe a que accede a las más profundas emociones humanas, igual que otras obras de arte como *2001: Una odisea del espacio*.

También resulta irónico que algunas voces hayan optado en tiempos recientes por alejar el campo de la IA de su mito feliz (o aterrador) en pos de un debate más «científico» —en otras palabras, han buscado descartar conceptos tan cargados emocionalmente como los de singularidad, consciencia e inteligencia sin dejar de beneficiarse de su presencia en la opinión pública—. Por ejemplo, es evidente que Russell desea separar la labor de corte serio en la inteligencia artificial general de los retratos que hace de ella la cultura popular en películas como *Ex Machina*. Russell considera que la consciencia es un problema filosófico ridículo (al fin y al cabo, lo ignoramos todo sobre ella), que los tests de Turing son ideas anticuadas y demasiado vulnerables al engaño, y que cualquier preocupación sobre la posibilidad de que las máquinas se pongan agresivas (o valientes, o emocionales de cualquier otra manera) resulta fundamentalmente errónea. Los ordenadores superinteligentes se limitarán a

perseguir sus objetivos. El problema es —y equivale a un riesgo existencial, incluso sin la imaginería propia de *Terminator*— que sus objetivos podrían no ser los mismos que los nuestros.

Russell admite que ese es un problema que ya tenemos con la IA. En concreto, llama la atención sobre los algoritmos de contenido-selección en la red, cuyo objetivo consiste en maximizar los ingresos publicitarios bombardeando a todo el mundo con anuncios pegajosos y relevantes. Es posible que a la IA superinteligente se le dé demasiado bien perseguir nuestros objetivos. Una metáfora apropiada es la del rey Midas, que obtuvo el poder de convertir todo lo que tocara en oro, y a quien de repente le resultó demasiado fácil convertir todo lo que tocaba en oro, incluyendo a su propia hija (ese no era el objetivo); de manera similar, la superinteligencia a la que se le encomendara un objetivo podría dar con una manera de alcanzarlo que acabara por eliminarnos, quizá incluso sirviéndose de los átomos de carbono de los propios seres humanos como un recurso más, como un medio más para alcanzar su fin.

Esa idea es una preocupación recurrente entre los aprensivos de la superinteligencia. Nick Bostrom plantea un escenario en el que la superinteligencia reciba la tarea en apariencia mundana de maximizar la producción de clips sujetapapeles (el objetivo que le han encomendado los seres humanos), y que poco a poco va convirtiendo todos los elementos del universo en una fábrica de clips sujetapapeles, incluyendo los elementos útiles de nuestro cuerpo. Eliezer Yudkowsky, exdirector del Instituto de Investigación de la Inteligencia de las Máquinas, afirmó una vez en broma: «La IA no te odia, tampoco te ama, pero es que estás hecho de unos átomos que ella podría usar para hacer otra cosa».¹⁷

La idea de que la superinteligencia venidera estará completamente centrada y será supercompetente en la consecución de un objetivo, pero a la vez mostrará un nulo sentido común, parece ir a contracorriente de lo que representa la superinteligencia misma —de la que se supone, al fin y al cabo, que es la suma de la inteligencia humana con algo más—. Científicos de IA como Gary Marcus, que entienden la inteligencia como algo dotado de (un mínimo de) sentido común (y quienes quizá también dispongan de él), han señalado que a un ordenador superinteligente que optimizara la fabricación de un producto humano de cara a su venta, como los clips sujetapapeles, también se le podría ocurrir que no debe destruir a los seres humanos, que son sus compradores. De nuevo encontramos una curiosa

simplificación de la superinteligencia implícita en las ansiedades de Russell y otros sobre la posibilidad de que esta actúe con un apego ciego e informático hacia el objetivo que le hayan encomendado sus programadores. Es extraño ponerse a vigilar algo así. El propio Russell admite que el sentido común y el lenguaje son hitos capitales que la IA no ha alcanzado aún. ¿A qué se debe su ausencia en esta imagen de la superinteligencia? Que los ordenadores del futuro dispusieran de sentido común permitiría obviar esas preocupaciones, a menos que después de todo acabaran revelándose como agresivos y diabólicos, cosa que a Russell le cuesta descartar como parte de este mito ridículo.

En todo caso, los escenarios del apocalipsis de los clips sujetapapeles sí que preocupan a algunos investigadores de mentalidad científica como Russell, quien sugiere que prevengamos esas posibilidades insertando en los ordenadores superinteligentes del futuro ciertos principios; el primero, asegurarse de que «no adjudican un valor intrínseco» a su propio bienestar, de modo que su único objetivo sea el de maximizar nuestras preferencias. El problema, tal y como nos recuerda Russell, es que a menudo no tenemos ni idea de cuáles son esas preferencias. En última instancia, somos propensos a no formular bien lo que deseamos, en el sentido de lo que le pasó al rey Midas.

Así, junto con el altruismo hacia la humanidad, la IA debe estar impregnada por el principio de la humildad, para desbaratar cualquier error que pueda cometer mientras intenta mostrarse altruista con nosotros (como el de convertir al director ejecutivo de la fábrica de clips sujetapapeles en un clip sujetapapeles llevada por la idea de que está exprimiendo de verdad, pero de verdad, hasta la última gota de su productividad al usar todos los medios que tiene a su alcance). Las máquinas altruistas y humildes ayudarán a protegernos del peligro de que los ejecutivos de empresas tecnológicas fumadores de puros (aunque lo más probable es que ya no los fumen) puedan conferirles motivaciones corruptas, y también de la posibilidad de que las máquinas se vuelvan demasiado listas de manera equivocada, y hagan algo parecido a transformar a todo el mundo en oro. Para mayor confusión, las máquinas «altruistas y humildes» también se acercan mucho al enfoque de *Ex Machina* sobre la IA: como algo «vivo» después de todo, con una inteligencia real (no solo dedicada a maximizar clips sujetapapeles) y una sensibilidad ética—. Se nos podría perdonar que

saquemos la conclusión de que el debate sobre la IA está condenado a mezclar la ciencia y el mito a perpetuidad.

Hay un tercer principio que Russell considera necesario para abortar la crisis existencial a la que nos conducirá la superinteligencia del futuro: la IA debería desarrollarse de manera tal que aprenda a predecir las preferencias de los seres humanos. En efecto, las máquinas deberían observarnos para aprender más acerca de lo que deseamos, lo que les permitirá descartar algunas acciones que podrían mandarlo todo al demonio, por así decirlo. Conocer las preferencias humanas permitirá que los ordenadores eviten hacernos daño mientras se esfuerzan por conseguir sus objetivos. (Russell no explica por qué deberíamos confiar en una superinteligencia que siga siendo lo bastante idiota para acabar con todos nosotros bajo la impresión de que nos está ayudando y otorgarle el papel de «panóptico» benevolente que no deje de observar y aprender nuestras preferencias.)

La versión de Russell de este relato sobre el riesgo existencial de la IA futurista nos trae a la cabeza los robots universales del dramaturgo checo Karel Čapek, que estaban diseñados para obtener una eficacia óptima en su trabajo y a los que se había despojado de manera deliberada de otros rasgos mentales como el sentido moral, la capacidad de apreciar la belleza y de experimentar sentimientos, y la consciencia. Por algún motivo, los autómatas supuestamente irreflexivos de su drama *R. U. R.* se sienten contrariados de todos modos y desencadenan una revolución robot que prácticamente acaba con la totalidad de la raza humana. Sin duda, el final de Čapek, es el motivo por el que recordamos esa obra de 1920. Nadie se emociona ante la perspectiva de que una Roomba tuneada aprenda de manera superinteligente (pero irreflexiva —ignora tú esa contradicción—) la mejor manera de aspirar la suciedad, o de limpiar la cocina, o de arreglar el coche. Sin duda resultaría de lo más útil, pero no es a eso a lo que nos referimos como superinteligencia. La que nos emociona es Ava. Una superinteligencia que no tenga consciencia ni sentimientos y que no sea capaz de cometer agresiones diabólicas no tiene nada de inteligente. Si carece de sentido común también será una pobre candidata para nuestras imaginaciones mitológicas. Será una calculadora.

Al unir las inteligencias humana y automática en lo que esencialmente es una búsqueda propia de la teoría de juegos para optimizar objetivos, Russell deja sitio para una visión en apariencia «científica» de la mente informática,

pero solo desde la limitación severa de las posibilidades de nuestra propia mente. Este vuelve a ser un error de inteligencia. La inteligencia humana es diversa, está llena aún de misterios profundos y, hasta donde sabemos, resulta eficazmente ilimitada. Al demoler la inteligencia humana, ligándola a una definición más sumisa ante la informática, el pensamiento contemporáneo sobre la IA tira por la borda una comprensión más rica de la mente. Y nos quedamos con un mundo simplificado.

Quizá ese mundo haga que la conversación sobre la llegada de la inteligencia artificial general parezca más razonable (porque «IAG» no equivale a gran cosa), pero a la vez pone en peligro el interés que existe en el proyecto mismo. Para el caso podríamos retirar la idea completa de superinteligencia y acometer un debate más sincero sobre la posibilidad muy real de que un virus informático con capacidad de destrucción mundial, pongamos que lanzado con evidente mala intención por un grupo de programadores, hiciera caer de manera irreflexiva los mercados financieros o accediera ilegalmente a unos datos capitales para la intimidad de los individuos o la seguridad de los países, y procediera a eliminar tales datos. Esa es una informática que se vuelve efectiva con el descubrimiento de una vulnerabilidad. Es el mundo real, no un mito.

EN RESUMEN

Podemos resumir estas posturas sobre la IA y la gente de la siguiente manera. Los kurzweilianos (mitólogos de la IA, punto) se entusiasman con la idea de que las máquinas posteriores a la singularidad vayan a disponer de conciencia, emociones, motivos y una vasta inteligencia. Con gran ironía, mantienen viva la exploración filosófica trasladando temas shakespearianos a la informática. (Los ordenadores vivirán ricas experiencias espirituales y serán grandes amantes, etc.) Podríamos denominar todo esto como el «efecto *Ex Machina*».

Los russellianos quieren que *Ex Machina* no sea más que una película, y reducir la conversación sobre la superinteligencia a ideas más respetables en lo matemático sobre una informática general que alcance «objetivos». Por desgracia, los russellianos también tienden a agrupar a los seres humanos según definiciones restringidas de la inteligencia. Eso reduce la

brecha que percibimos entre humanos y máquinas, pero siempre a cambio de limitar las posibilidades de los primeros. Los russellianos son líderes de pensamiento en una tendencia cultural que he bautizado como «el mundo simplificado». Tal y como dice Jaron Lanier, «una nueva generación ha alcanzado la mayoría de edad con unas expectativas reducidas de lo que puede ser cada persona, y sobre aquello en lo que se puede convertir».¹⁸

Tanto los kurzweilianos como los russellianos proclaman una visión tecnocéntrica del mundo que a la vez simplifica la visión de la gente —en especial, con sus enfoques deflacionarios sobre la inteligencia como una forma de computación— y expande las visiones de la tecnología al promover un futurismo dedicado a la IA como ciencia y no como mito.

Centrarnos en el batitraje y no en Bruce Wayne nos ha metido en un montón de problemas. Vemos posibilidades ilimitadas para las máquinas, pero un horizonte restringido para nosotros. De hecho, la inteligencia futura de las máquinas es una cuestión científica, no mitológica. Si la IA continúa siguiendo el mismo patrón de superar las expectativas en el mundo falso de los juegos o la colocación de anuncios podríamos acabar, al límite, con unos sabios increíblemente indiscretos y peligrosamente idiotas.

Ahora vamos a centrarnos en la ciencia de la IA, y es aquí —en la investigación científica misma— donde el mundo simplificado se vuelve complejo (y misterioso) de nuevo. Y es que, cuando nos liberamos de las ataduras de nuestros errores de inteligencia, se nos cae la venda de los ojos y un problema verdaderamente formidable se presenta ante nosotros.

Segunda parte
EL PROBLEMA DE LA
INFERENCIA

Capítulo 8

No calcules, analiza

La IA es la búsqueda de la inteligencia. A lo largo de los varios capítulos que conforman esta parte del libro espero convencer de que esta búsqueda se enfrenta a obstáculos importantes, que no sabemos cómo superar. Para ello, debemos investigar la naturaleza misma de la inteligencia. Y no hay nada mejor para iniciar nuestra investigación que acompañarnos de un «joven extraño e interesante», el detective aficionado Auguste Dupin, a quien nos presenta un narrador sin nombre en el que quizá sea el primer relato de detectives de la historia, «Los crímenes de la calle Morgue».¹

SOBRE LA RESOLUCIÓN DE CRÍMENES

El narrador —que comparte numerosos rasgos con Edgar Allan Poe, el autor— nos cuenta desde un principio que está obsesionado con los métodos de pensamiento. Siente curiosidad por la manera en que la mente humana conecta fragmentos de información que en apariencia no tienen relación entre sí a partir de la observación cuidadosa y el razonamiento —a partir de inferencias—. Menuda casualidad, pues, que el narrador se encuentre alojado en una vieja casa con Dupin, y que pueda pasarse todo el día cerca de tan brillante detective.

No tardamos en averiguar que Dupin no es un sujeto normal. Posee ese tipo de personalidad extraña de quienes son verdaderamente originales. Y vaya si es extraño. Dupin procede de una familia distinguida, pero se ha visto rebajado a un estado cercano a la pobreza, cosa a la que apenas presta atención porque está pensando constantemente, perdido en sus ideas. Cuando habla se debe a que está rumiando en voz alta. Por supuesto, eso podría acabar resultando molesto. Pero el narrador siente gran aprecio por su «peculiar capacidad analítica». Dice: «Nos pasábamos el día leyendo, escribiendo o conversando, hasta que el reloj nos advertía de la llegada de la Oscuridad verdadera. Entonces salíamos a las calles cogidos del brazo, seguíamos hablando sobre los temas del día o vagábamos de aquí para allá hasta una hora tardía, en busca, entre las luces y sombras salvajes de la ciudad atestada, de la excitación mental sin fin que la observación callada suele procurar».

Dupin es un prototipo, un modelo, como Sherlock Holmes. E, igual que Holmes, repara en lo que la policía se las arregla para pasar por alto en la aplicación de su «simple diligencia y actividad».

Una noche, a solas en la vieja casa con Dupin, el narrador coge la edición vespertina de la *Gazette des tribunaux* y se entera de los asesinatos de la calle Morgue:

Crímenes extraordinarios. — La madrugada pasada, alrededor de las tres, los vecinos del Quartier St. Roch despertaron de su sueño a causa de una sucesión de alaridos terribles, procedentes, al parecer, del cuarto piso de una casa de la calle Morgue de la que solo constaban como ocupantes una tal Madame L’Espanaye y su hija, Mademoiselle Camille L’Espanaye. Con cierto retraso, debido a un intento infructuoso por entrar de la manera tradicional, se utilizó una palanca para romper la puerta y entre ocho y diez vecinos accedieron al lugar acompañados de dos gendarmes. En aquel momento los gritos ya habían cesado pero, mientras el grupo subía a la carrera el primer tramo de escaleras, procedentes al parecer de la parte superior del edificio, se pudieron distinguir dos o más voces roncadas, que mantenían una discusión furiosa. Al llegar al segundo descansillo, los sonidos también habían cesado y reinaba un silencio perfecto. El grupo se dispersó, sus miembros se apresuraron a entrar en una habitación tras otra. Al llegar a la espaciosa sala trasera del cuarto piso (la puerta de la cual hubo que forzar, al estar cerrada con llave y encontrarse la llave del lado de dentro), el espectáculo que se desplegó ante los presentes sumió a cada uno de ellos no tanto en el horror como en la perplejidad.

El apartamento se encontraba completamente revuelto... los muebles estaban rotos, los habían arrojado en todas direcciones. Solo había una armadura de cama, a la que le habían quitado el colchón para arrojarlo en medio del suelo. Sobre una silla descansaba una navaja, manchada de sangre. En la chimenea había dos o tres mechones largos, gruesos y canosos de cabello humano, también salpicados de sangre y que parecían haber sido arrancados de raíz. En el suelo se encontraron cuatro Napoleones, un pendiente de topacio, tres cucharas grandes de plata, otras tres más pequeñas de argén y dos bolsas que contenían cerca de cuatro mil

francos en oro. Los cajones del secreter, en una esquina, estaban abiertos y al parecer los habían registrado, aunque en ellos quedaban numerosos artículos. Debajo de la cama (no debajo del armazón de la cama) encontraron una pequeña caja fuerte de hierro. Estaba abierta, la llave aún en la puerta, y no contenía nada más que algunas cartas viejas y otros documentos de escasa importancia.

No había rastro de Madame L'Espanaye pero, al reparar en la inusual cantidad de hollín que se amontonaba en el hogar, se realizó una búsqueda en la chimenea y (¡cuán horrible es contar esto!) de su interior sacaron a rastras el cadáver de la hija, que colgaba boca abajo. Lo habían introducido a la fuerza por la estrecha abertura y había recorrido una distancia considerable. El cuerpo aún conservaba buena parte de su calor. Tras examinarlo se reparó en que había sufrido diversas escoriaciones, sin duda debidas a la violencia con que la habían metido chimenea arriba y después retirado. Tenía varios rasguños fuertes en la cara y unos cardenales oscuros en la garganta, y las marcas profundas de unas uñas, como si la fallecida hubiera sido asfixiada hasta la muerte...

Tras una investigación minuciosa de cada parte de la casa, al no descubrirse nada más, el grupo se dirigió hacia un pequeño jardín empedrado en la parte trasera del edificio, donde yacía el cadáver de la anciana, que había sufrido un corte tan profundo en la garganta que, en el intento por levantarla, la cabeza rodó por el suelo. El cuerpo, igual que la cabeza, había sido terriblemente mutilado —el primero, a un extremo tal que a duras penas conservaba apariencia alguna de humanidad.

De momento, entendemos que no existe la más ligera pista sobre este misterio horrible.²

Al día siguiente, la *Gazette* publica nuevos detalles sobre el caso. A partir de los recuentos de quienes prestan declaración podemos reunir la información relevante. Madre e hija eran pudientes. Tres días antes de los asesinatos, la madre había retirado una suma importante de dinero del banco, en oro, que tras los crímenes se encontró a plena vista, sin tocar, en el suelo. También resulta curioso que uno de los policías del primer grupo que llega a la escena manifieste haber oído dos voces; una, de manera evidente, de un hombre que hablaba francés, y otra que no pudo reconocer en absoluto y que definió como «severa, aguda y muy extraña». Pensó que se trataría de un extranjero, posiblemente español. Más adelante, otros testigos describirían la voz como posiblemente italiana, rusa o inglesa.

Era desconcertante. El dinero —quizá el motivo más probable para cometer un asesinato— se queda sin tocar en la casa. Las puertas están cerradas con llave por dentro. El cuerpo de la hija se encuentra en el interior de la chimenea, alojado en ella con tanta fuerza que es necesaria más de una persona para sacarlo de allí. Y las voces también. Las de los asesinos, al parecer, pero, aunque la policía oyó con claridad, mientras subían las escaleras de la casa, que eran dos, solo una resulta reconocible, y afirman que la otra es una mezcla extraña de jergonza. Ninguno de los testigos es

capaz de decir con exactitud lo que se dijo (si se dijo algo), ni en qué idioma.

La policía está perpleja. Las declaraciones de los testigos no hacen más que ahondar su confusión. Con ello vengo a decir que todas las pistas, en su conjunto, no apuntan hacia ninguna parte. Los crímenes son un misterio, y ese es el motivo por el que nuestro excéntrico detective aficionado, Dupin, siente un interés tan entusiasta.

El narrador sugiere que Dupin resuelve el caso muy pronto, tras leer en el periódico el comunicado de la policía. No obstante, los dos reciben permiso para visitar la vieja casa de la calle Morgue con la escena del crimen aún intacta. De regreso a casa, Dupin se detiene en las oficinas de otro periódico y pone un anuncio en la sección de objetos perdidos. ¿Alguien en París, presumiblemente el marinero de un navío maltés, ha perdido a su orangután? Existe la posibilidad de que el dueño llame para reclamarlo.

Y aquí aparece la inferencia que permite resolver el caso de los crímenes de la calle Morgue: ningún ser humano mató a la anciana y a su hija aquella noche. El asesino no fue una persona, sino un animal salvaje al que un marinero se trajo de la jungla y mantuvo en alguna morada cercana. Frenético, habiendo escapado de su amo, el orangután, tras balancearse en el postigo externo, entró de un salto por la ventana de la casa profiriendo chillidos y alaridos y esgrimiendo una navaja de afeitar. Ahí está el arma del crimen: en la navaja que decapita a la anciana y en la pura fuerza bruta del animal, que mete a la hija en la chimenea arrastrándola por los pies.

¿Y la voz humana que oyeron los testigos? El dueño del orangután. ¿Y los sonidos amortiguados e ininteligibles? Los gruñidos del animal.

EL MÉTODO CONJETURAL

Pero ¿quién inferiría eso a partir de los datos del caso? Sin duda, todo el mundo los tiene delante de los ojos. En realidad, Dupin se limitó a conjeturar. La policía siguió los métodos tradicionales hasta que estos no desembocaron en nada. A partir de ese momento también comenzaron a hacer conjeturas. La única diferencia radicó en que la de Dupin fue la buena.

Poe comienza «Los crímenes de la calle Morgue» reflexionando sobre la naturaleza del pensamiento. La historia ficticia del crimen se inicia desde la no ficción. El autor busca las palabras adecuadas. El razonamiento de Dupin, decide, es un triunfo del análisis, en contraste con cálculos más formalistas. El cálculo consiste en conectar unos puntos que ya se conocen; en aplicar las reglas del álgebra, pongamos. El análisis consiste en ofrecer sentido a esos puntos, dando saltos o efectuando suposiciones que los expliquen —y, a continuación, desde una cierta perspectiva, usando el cálculo para comprobarlas—. El cálculo tiene sus límites: «Pero es en aquellas cuestiones que trascienden los límites de la mera regla cuando se demuestra la habilidad del analista». Seguir las reglas no es suficiente, pero no queda claro qué más hace falta con exactitud. El aprecio de Poe hacia ese misterio se vuelve evidente en la declaración que realiza al comienzo de la historia: «Las condiciones mentales que suelen considerarse analíticas son, en sí mismas, poco susceptibles al análisis».³

Algunas décadas más tarde, el científico y filósofo norteamericano Charles Sanders Peirce iba a leer los relatos de Poe con gran fascinación. Peirce también se preguntó por la manera en que pensamos, en que razonamos acerca de las cosas. Incluso se las arregló para capturar la gimnasia mental de Dupin en forma de símbolos lógicos. No sabía cómo automatizar el perspicaz estilo conjetural del detective, pero pensó que se trataba de un aspecto básico del pensamiento humano en general.

Para Peirce, el pensamiento no es un cálculo sino un salto, una suposición. No hay ninguna certeza. Nos dedicamos a juntar las piezas. Explicamos y revisamos. Al haber vivido durante el siglo XIX, Peirce no llegó a conocer los ordenadores digitales, pero anticipó lo que iba a hacer de la IA un problema tan complicado para todo el mundo. En realidad, se reduce a esto: puesto que nuestro propio pensamiento se basa en una serie desconcertante de conjeturas, ¿qué esperanza hay de que podamos programarlas?

Con el tiempo, Peirce acabaría desarrollando todo un marco explicativo para el razonamiento humano que se basó en la lógica formal y sus clases, como la deducción y la inducción.

Pero había un tercer elemento, reflexionó Peirce, que capturaba nuestros juegos conjeturales. Lo llamó «abducción», y en él nos vamos a centrar a continuación.

Capítulo 9

El puzle de Peirce (y el rompecabezas de Peirce)

Para aquellos que estén familiarizados con su trabajo, Charles Sanders Peirce forma parte de un selecto grupo de pensadores importantes y verdaderamente originales: En la biografía *Charles Sanders Peirce: A Life* [«Charles Sanders Peirce. Una vida»], el historiador Joseph Brent se refiere a él como, «quizá, la mente más importante que Estados Unidos haya producido nunca». En 1934, el filósofo Paul Weiss describe a Peirce en *The Dictionary of American Biography* [«Diccionario de biografías norteamericanas»] como «el más original y versátil de los filósofos norteamericanos, y el más importante lógico del país». El crítico e historiador cultural Lewis Mumford le situó en compañía de genios iconoclastas como Roger Bacon y Leonardo da Vinci. Y, durante una entrevista de 1976, cuando le preguntaron al pionero de la ciencia lingüística del MIT Noam Chomsky por sus influencias, este contestó: «En relación con las cuestiones que hemos estado tratando [relativas a la filosofía del lenguaje], el filósofo que siento más cercano y aquel al que prácticamente parafraseo es Charles Sanders Peirce».¹

BRILLANTE PERO ABANDONADO

Como Albert Einstein, Peirce era zurdo y pensaba en imágenes. Bosquejaba inferencias lógicas en forma de diagramas. Pasó los últimos años de su vida escribiendo solo en casa, quejándose de que pasaba hambre y frío, y de que era demasiado pobre para permitirse el combustible del fogón. Los pocos amigos que tenía se preocupaban por él y se las apañaron para conseguirle una serie de conferencias en Harvard sobre los fundamentos de la lógica, en las que Peirce explicó los tipos de inferencia lógica desde un marco que, según él, afianzaba el método científico —un método para pensar con claridad—. Entre los asistentes se encontraba William James, el famoso filósofo y precursor de la psicología de Harvard, quien más tarde confesó que no había acabado de entender las conferencias, que las matemáticas relacionadas con las imágenes y diagramas de Peirce se encontraban más allá de su comprensión. Al parecer, James no fue el único a quien le sucedió algo así; las conferencias pasaron mayormente desapercibidas y solo aparecieron en forma de libro décadas después.

Nacido en el entorno de la cultura científica victoriana de Cambridge, Massachusetts, en 1839, Peirce fue miembro de una familia pudiente y destacada en diferentes ámbitos. Su padre fue un eminente profesor de matemáticas en Harvard. Un primo más joven, Henry Cabot Lodge, se convirtió en un poderoso senador. Peirce recibió una educación clásica y en 1863 se graduó *summa cum laude* en la Escuela Científica Lawrence de la universidad de Harvard. Pasó treinta años trabajando como investigador científico en el Servicio de Costas y en el Servicio Geodésico Nacional, donde realizaba estudios topológicos de la superficie terrestre usando precisas medidas gravimétricas. Fue químico aficionado, un prestigioso conferenciante sobre lógica en la universidad Johns Hopkins y el primer delegado norteamericano de una asociación científica internacional. Fue científico, lógico, filósofo, escritor, prolífico reseñista de libros para el *Nation*, y más. Max H. Fisch, estudioso de Peirce que se pasó décadas investigando su vida y obra, ofrece este adecuado juicio grandilocuente sobre sus numerosos logros:

¿A quién pertenece el intelecto más original y versátil que hayan generado hasta el momento las Américas? La respuesta «Charles S. Peirce» no encuentra oposición, porque cualquier nombre que citáramos en segundo puesto se hallaría tan lejos que no valdría la pena ni mencionarlo. Matemático, astrónomo, químico, geodesta, topógrafo, cartógrafo, metrólogo, espectrólogo, ingeniero, inventor; psicólogo, filólogo, lexicógrafo, historiador de la ciencia, economista matemático, estudioso de la medicina a lo largo de toda su vida; crítico de libros, dramaturgo, actor, autor de relatos; fenomenólogo, semiótico, lógico, retórico, metafísico... Es el único filósofo de sistemas de las Américas que se ha mostrado a la vez competente y productivo en

lógica, matemáticas y en una amplia gama de ciencias. Si en ese sentido ha tenido algún par a lo largo de toda la historia de la filosofía, no habrán sido más de dos.²

Sin embargo, Peirce murió como un paria, en gran medida olvidado. Los genios olvidados han sido tan comunes a lo largo de la historia que a veces los acabamos redescubriendo, como sucedió con Tesla. Pero, sin duda más que Tesla —quien, al fin y al cabo, alcanzó una especie de fama póstuma como motivo de inspiración de Elon Musk para bautizar una compañía de coches eléctricos en su honor—, Peirce sigue siendo un pensador importante al que los libros de historia no suelen mencionar. Su trabajo ha sido apreciado sobre todo en la filosofía, donde se le recuerda como el fundador de la escuela conocida como pragmatismo.

Sus obras tempranas sobre computación están completamente olvidadas. Los académicos continúan sondando sus textos voluminosos sobre la naturaleza de la lógica, pero se trata de un tema oscuro, demasiado difícil para conectarlo con el debate general. Por ello, aunque quienes han comprendido la importancia y el relieve de sus ideas sobre la naturaleza de la lógica le hayan comparado con Aristóteles, el debate contemporáneo sobre las ideas de Peirce requiere en la mayoría de los círculos un bosquejo biográfico... y una explicación, incluso una disculpa.

FÍSICA, FILOSOFÍA Y PERSONALIDAD

¿Por qué se ha olvidado a Peirce? Su vida personal nos ofrece una pista: logró irritar a casi todo el mundo. William James no dejó de ser su amigo íntimo de por vida, pero incluso él salió de su primer encuentro con Peirce, cuando ambos eran alumnos de Harvard, con sensaciones encontradas, tal y como le contó por carta a su familia: «Está el hijo del profesor Peirce, de quien sospecho que es un “tipo” muy inteligente y con mucho carácter, bastante independiente pero también agresivo».³ Comprensivo, James más tarde se refirió a Peirce como «ese ser extraño y rebelde».⁴

La personalidad espinosa de Peirce y su indiferencia hacia las costumbres de sus contemporáneos le granjearon problemas sin fin, tanto en lo personal como en lo profesional. Solía ofender a los miembros de la alta sociedad victoriana de Nueva Inglaterra (incluyendo a su familia), quienes le rehuían

justificadamente. Harvard se negó a ofrecerle una cátedra a causa de una conocida infidelidad en su matrimonio. El Servicio de Costas de Estados Unidos, para el que trabajó durante décadas, acabó despidiéndole por no entregar sus informes a tiempo y por perder un equipo de gran valor mientras viajaba por Europa. También le despidieron de la Johns Hopkins tras unas denuncias indeterminadas relacionadas con una conducta impropia. Hoy diríamos que no encajaba en los sitios, que fue el estereotipo perfecto del genio incomprendido. Su constitución era incapaz de respetar las reglas del juego.⁵

Los escándalos e idiosincrasias personales de Peirce ayudan a explicar que, por ejemplo, los registros de su vida privada —volúmenes de documentos— permanecieran sellados en la biblioteca Houghton de Harvard hasta 1956, cuarenta y dos años después de su muerte. Sus artículos científicos y filosóficos —incluyendo muchos de un interés enorme para la ciencia computacional y, en especial, la IA— se quedaron relativamente intactos, e inéditos, en los archivos de Harvard, a falta de «unos pocos miles de dólares» como «garantía por los gastos iniciales de su publicación», en palabras de Lewis Mumford.⁶ Al no entender sus ideas y su valor, por su deseo de no atraer el escándalo, muchos de quienes conocieron a Peirce jamás fueron testigos del restablecimiento de su reputación, ni de la publicación de buena parte de la obra de su vida. Peirce mismo murió en el olvido; le sobrevivió Juliette, su no menos enigmática segunda esposa, una francesa que también había llevado una vida accidentada. De manera apropiada, la familia de Peirce a veces describía a Juliette como una paria, o una «gitana».

Solo mucho, mucho tiempo después hemos llegado a comprender la enormidad de la contribución de Peirce a las matemáticas y, en especial, a la lógica —hasta cierto punto, ni siquiera hemos acabado de hacerlo—. Son de la máxima importancia sus ideas sobre la inferencia lógica y, en particular, para desplazarnos hacia la atracción principal de su trabajo vital, su exploración de la profundidad y el misterio de lo que él denominó la «inferencia abductiva», una especie de conjetura explicativa que, según se dio cuenta, afianza la mayor parte de nuestro pensamiento.

Peirce reparó en el hecho de que el razonamiento abductivo había quedado fuera de todos los recuentos de razonamiento lógico desde tiempos de Aristóteles. Tampoco encajaba en el marco lógico habitual que asumían las matemáticas y los cursos de lógica. Vio la abducción como una pieza

lógica perdida que planteaba preguntas fundamentales sobre la automatización y la inteligencia. Si hubiera sabido de la IA, lo más probable es que hubiera visto algo que todavía hoy se sigue pasando por alto a menudo: que el problema de la inferencia abductiva confronta a la IA con su desafío central, aún no resuelto del todo.

EL PUZLE DE LA INFERENCIA

El narrador de Edgar Allan Poe se devanó los sesos buscando las palabras que describieran aquello a lo que Peirce más tarde dedicaría diversos volúmenes: la inferencia abductiva. Pero la inferencia abductiva es un tipo de inferencia. ¿Qué es una inferencia? Un nombre, para comenzar. La forma verbal es «inferir», que se refiere a una acción. Etimológicamente, «inferir» significa «traer aparejado», del latín «*in*», en, y «*ferre*», traer. El *Oxford English Dictionary* nos dice que se trata de algo que hacemos cognitivamente, con la mente: «Alcanzar una opinión o decidir que algo es verdadero basándose en la información de la que se dispone».

Por desgracia, el *OED* también dice que «deducir» es sinónimo de inferir, lo cual no nos ayuda demasiado (porque la deducción es solo una de las formas de la inferencia).

El *OED* también ofrece algunos ejemplos de uso que recalcan la generalidad de la palabra «inferir» en el lenguaje cotidiano.

Inferir algo (a partir de algo): Hay que inferir buena parte del sentido a partir del contexto.
Inferir los motivos del asesino queda en manos de los lectores.
Inferir que: Resulta razonable inferir que el gobierno conocía esos acuerdos.

La inferencia consiste en dar lugar a una nueva idea, lo cual en lógica equivale a extraer una conclusión y, de manera más general, implica usar aquello que ya sabemos, y lo que vemos u observamos, para actualizar convicciones previas. Podríamos inferir los motivos del asesino (por tomar prestado uno de los ejemplos del *OED*) sirviéndonos de lo que ya sabemos y de lo que hemos leído en los periódicos (igual que Dupin).

La inferencia también es una especie de salto, pero un salto que estimamos razonable, como cuando inferimos que «el gobierno conocía esos acuerdos» —de nuevo, de acuerdo con los conocimientos previos de

los que dispusiéramos (conocimiento público o compartido) junto con la lectura (observación) de alguna historia o historias de actualidad.

La inferencia es un acto cognitivo básico de la mente inteligente. Cuando un agente cognitivo (una persona, un sistema de IA) no es inteligente, sus inferencias serán incorrectas. Pero todo sistema que infiera algo debe disponer de una inteligencia básica, porque el acto mismo de utilizar lo que se sabe y lo que se observa para actualizar nuestras convicciones se encuentra ligado de manera inevitable a lo que queremos transmitir al hablar de inteligencia. El sistema de IA que no infiere nada no merece recibir el nombre de IA. (Aunque podríamos decir que incluso un sistema encargado de etiquetar fotos de gatos infiere que lo que «ve» es un gato, así que el listón puede encontrarse bastante bajo.)

Es imposible pillar un chiste, descubrir una nueva vacuna, resolver un asesinato como hace Dupin o simplemente mantenerse al día de los diversos acontecimientos y comunicaciones que se dan en el mundo sin algún tipo de capacidad de inferencia. Sabemos un montón de cosas, sin duda, pero solo las inferencias nos permiten obtener conocimientos (o convicciones) nuevos. Sabemos que mañana saldrá el sol, así que no necesitamos inferirlo. Asimismo, no nos molestamos en inferir que seguimos teniendo la mano pegada al brazo. Es un conocimiento del que ya disponemos, un conjunto de convicciones que ya hemos formado. Pero nuestro conocimiento cambia de manera constante, se actualiza. Si vemos por la ventana que ha oscurecido misteriosamente, pues es demasiado temprano para ello, podemos inferir que se está produciendo un eclipse solar, o quizá que una enorme tormenta de arena ha ocultado el sol por el oeste, o quizá que ha habido un holocausto nuclear. Todo depende de... ¿qué sabemos en ese momento? Por lo que vemos, ¿qué es lo que tiene más sentido?

En un sentido amplio, siempre estamos infiriendo; es como una de las condiciones de estar despierto. Puedo ir a la cocina, encontrarme con una lata medio vacía de Pepsi e inferir que mi hermana la ha dejado allí, ya que bebe Pepsi y está de visita. Por otro lado, hay unos operarios reformando las encimeras, y antes reparé en que uno de ellos se bebía una Pepsi. Es más, hace un rato estuve bebiéndome una Pepsi y la dejé sin terminar en el porche, así que es posible que mi esposa la haya entrado a la casa. Acabamos conjeturando una explicación que tenga sentido, dado lo que sabemos y el contexto en el que nos encontramos. Se trata de una inferencia «en tiempo real», ya que extraemos esas conclusiones mientras entramos en

la cocina. Las circunstancias del mundo real no dejan de cambiar, así que la inferencia en tiempo real es normal. Al fin y al cabo, pensamos dentro del tiempo. Un programa informático que tarde diez mil millones de años en resolver un problema no será para nada inteligente, y tampoco lo será el que, en tiempo real, se estrelle contra una pared.

La naturaleza provisional de numerosas inferencias implica que las iniciales puedan estar equivocadas, sobre todo cuando nos hemos apresurado a alcanzarlas. Si llego tarde a la oficina, el jefe puede inferir que no me tomo las cosas con seriedad, cuando de hecho me he encontrado con un atasco de tráfico a causa de un accidente. En otras palabras, el jefe ha extraído una conclusión basada en una impresión preformada o prejuicio sobre mí. La gente usa la palabra «inferencia» en sus conversaciones cotidianas otorgándole ese sentido, refiriéndose a un salto apresurado hacia una conclusión injustificada: «Oh, es Suzy, que está infiriendo todo tipo de locuras sobre ti después de lo que dijiste anoche». Y es cierto que, en un sentido técnico, Suzy está realizando inferencias, pero la implicación aquí es que se trata de inferencias tendenciosas, y que Suzy se encuentra demasiado dispuesta a hacer suposiciones injustas (quizá porque está de mal humor o porque no le caes bien).

En un sentido más específico, la inferencia ingresó en el léxico matemático hace mucho tiempo, y en épocas más recientes ha participado en los debates sobre la informática y la IA. En ese contexto, la «inferencia en tiempo real» puede referirse a un robot que se oriente por un entorno dinámico, como una calle bulliciosa. La «inferencia probabilística» extrae conclusiones de datos estadísticos y tiene una aplicación evidente en los enfoques de IA basados en datos.

Hubo un tiempo en que los científicos de IA se pelearon vigorosamente con una condición previa de la inferencia, cuyo uso inteligente ya conocemos: la cuestión del «conocimiento». Los sistemas que no saben nada tampoco pueden inferir gran cosa. Así que aquellos primeros investigadores intentaron programar el conocimiento en los sistemas de IA, para ayudarlos a que dieran sentido a su sensor o a sus impulsos de texto. Se descubrió (por las malas, a través de fracasos repetidos) que los sistemas de IA con amplios depósitos de conocimiento en forma de datos y reglas seguían necesitando usar ese conocimiento dentro de un contexto para extraer conclusiones relevantes. Ese «uso» del conocimiento es lo que lleva a que las inferencias sean tan difíciles. ¿Qué fragmento de conocimiento

resulta relevante en el pajar de mi memoria para aplicarlo al mundo dinámico y cambiante que me rodea?

La capacidad de determinar los fragmentos de conocimiento con relevancia no es una aptitud informática. Poe insiste en que, en el reino de lo «analítico», no existen fórmulas que permitan llegar a la percepción humana; se trata de «asuntos que van más allá de los límites de la mera regla» o cálculo. En efecto, Dupin parece llegar a una explicación para los crímenes —el orangután— a través de una conjetura casual, que más tarde verifica al encontrarse con el dueño del animal perdido. Por tanto, ¿se trató solo de una suposición por su parte? En un sentido importante, sí. Pero eso no anula su carácter de inferencia. Hace que esta cobre importancia.

MÁS TURING

En su artículo seminal de 1950, «Maquinaria computacional e inteligencia», Turing desestimó la cuestión de que las máquinas llegaran a pensar burlándose de su propio título y asegurando que el «pensamiento» es subjetivo y carece por completo de rigor científico. Decir que los ordenadores piensan es como decir que los submarinos nadan. Que se hable de «nadar» ya es un ejercicio de antropomorfismo. Los delfines nadan, pero los submarinos no. Turing pensó que con el uso de la palabra «pensamiento» sucedía algo parecido. Cuando un ordenador juega al ajedrez, ¿quién puede decir si está pensando o si se limita a hacer cálculos?

A Turing le interesaba una mente completamente programable. Por consiguiente, desechó la distinción que había realizado en un principio entre intuición e ingenio llevando la intuición —fuera lo que fuese— a la esfera de la computación. De ese modo, logró que el problema de la IA fuera totalmente comprobable. La tesis era radical incluso según sus propios parámetros previos, pero no se lo tendremos en cuenta porque sentó las bases para que, más tarde, en esa misma década, los investigadores de IA se pusieran a trabajar sin preocupaciones filosóficas que retrasaran sus avances.

Por desgracia, lo que nunca se abordó debidamente fue la manera exacta en que la inferencia informática podía ser igual que —o podía convertirse en— una inferencia humana. La disciplina no comenzó con una teoría de la

inferencia, lo que podría haber proporcionado un plano completo para su desarrollo futuro (o la demostración de su imposibilidad). Que los investigadores de IA carezcan de una teoría de la inferencia es como si los ingenieros nucleares se hubieran puesto a trabajar en la bomba nuclear sin haber resuelto antes los detalles de las reacciones de fisión. Es evidente que no basta con saberse la ecuación de Einstein. Y que tampoco basta con que los entusiastas de la IA tengan conocimientos de teoría de la computación, porque la pregunta misma a la que se enfrentan los científicos que trabajan con la IA es cómo se puede llevar la computación al rango adecuado y a los tipos de inferencia que exhibe la mente. Era algo que se tendría que haber preguntado desde un principio. Al ignorar la pregunta o esquivarla, la disciplina generó falsas esperanzas, condujo a callejones sin salida y perdió el tiempo de manera inevitable.

Y es que hay mucho sobre lo que reflexionar. Tomemos, por ejemplo, las numerosas inferencias que encontramos en la historia de la ciencia. Los científicos bosquejan hipótesis, y a continuación las ponen a prueba. Pero a esas hipótesis no se ha llegado de manera mecánica; son famosas por brotar en la cabeza de los científicos (por lo general, después de que estos hayan alcanzado la maestría en su disciplina). Tal y como hizo Turing una vez, los estudiantes de esos descubrimientos científicos tienden a apartar los saltos intelectuales de las formalidades de la práctica científica, de modo que el acto central de inteligencia «se suma al viaje sin pagar por él» —y se queda sin que lo sometan a análisis—. Pero esas hipótesis son actos genuinos de la mente, primordiales para cualquier ciencia, y a menudo no se los puede explicar señalando los datos o las pruebas o cualquier cosa evidente o programable.⁷

Al plantear que la Tierra giraba alrededor del Sol, y no al revés, Copérnico ignoró montañas de datos y pruebas acumuladas a lo largo de los siglos por los astrónomos que habían trabajado con el antiguo modelo ptolemaico. Lo redibujó todo, con el Sol en el centro, y calculó un modelo heliocéntrico que fuera útil. Aún más importante, el modelo copernicano resultó de hecho menos predictivo pese a ser correcto. En un principio se trató solo de un armazón que, en caso de completarse, podría ofrecer una serie de explicaciones elegantes para reemplazar otras cada vez más enrevesadas, como la aparente retrogradación de los planetas, que era una pesadilla para el modelo ptolemaico. Solo ignorando en un principio todos los datos, o reconceptualizándolos, pudo Copérnico rechazar el modelo

geocéntrico e inferir una estructura radicalmente nueva para el Sistema Solar. (Y fíjate en que eso plantea la pregunta de cómo le habrían podido ayudar los big data, dado que todos los datos encajaban con el modelo equivocado.)

El salto copernicano que hizo despegar la revolución científica podría describirse mejor como una suposición inspirada. Y lo mismo podría decirse acerca de la elección por parte de Kepler de la elipse para describir el movimiento de los planetas, porque en las órbitas planetarias se puede encajar un número muy grande (técnicamente infinito) de figuras geométricas (quizá excluyendo las de tipo trascendental, como las sinusoides). La elipse no era una solución más sencilla que las demás —no se trató de una explicación del tipo navaja de Occam—. Kepler literalmente conjeturó una explicación que a él le «parecía correcta». Que las conjeturas desemboquen en descubrimientos no encaja en el relato mecánico de la ciencia; más bien lo contradice. Pero la labor del detective, los descubrimientos científicos, la innovación y el sentido común son todos obra de la mente; son todos inferencias que los científicos de IA que buscan máquinas de inteligencia general deben explicar de alguna manera.

Como se puede ver, la creación de modelos cognitivos —la construcción de un ordenador que piense, que infiera— resulta desconcertante. Los investigadores en IA (al menos de momento) deberían preocuparse sobre todo de la inferencia en su contexto cotidiano. ¿Por qué? Pues porque la vasta mayoría de las inferencias que realizamos son aparentemente mundanas, como las suposiciones y saltos variopintos que efectuamos en el transcurso de una conversación común y corriente. Por desgracia para los investigadores en IA, ni siquiera esas inferencias mundanas son fáciles de programar. El test de Turing, por ejemplo, resulta complicado en esencia porque la comprensión del lenguaje natural requiere un montón de inferencias de sentido común, que ni son ciertas lógicamente ni resultan (a menudo) demasiado probables. Requiere, en otras palabras, muchísimas abducciones.

Por lo general ni siquiera reparamos en esas inferencias, lo cual está bien: en caso de hacerlo, tenderíamos a quedarnos atascados en bucles solipsistas, dándoles vueltas a las cosas en nuestras cabezas. Esto nos devuelve a Peirce y, de manera más específica, nos conduce al marco tripartito de inferencias en el que se afianza la inteligencia: deducción, inducción y abducción.

Capítulo 10

Problemas de deducción e inducción

A lo largo de la mayor parte de la historia intelectual, la inferencia ha sido sinónimo de deducción. Aristóteles estudió una forma simple de la deducción conocida como el silogismo —dos proposiciones de las que se sabe o se cree que son ciertas conducen a una tercera, la conclusión—. Aristóteles desarrolló una forma temprana de lógica usando silogismos para analizar los argumentos que realizaban él y otros, y para sentar las bases del razonamiento correcto. En su tradición, la inteligencia debe cumplir con las reglas conocidas de la deducción.

Eso tiene sentido. No deberíamos dejarnos persuadir, por ejemplo, por una persona que argumentara que Ray Charles es Dios porque «Dios es amor» y «El amor es ciego» (igual que Ray Charles). Ese argumento es falaz; quebranta las reglas del razonamiento deductivo. Precisamente, anotar todo eso forma parte de la tradición de la lógica deductiva. Aristóteles también exploró la relación entre las reglas deductivas y el llamado razonamiento práctico —por ejemplo, cuando un agente inteligente formula un plan para alcanzar un objetivo cuyos pasos se pueden analizar lógicamente—. (El plan podría ser «correcto» de forma comprobable y, sin embargo, fracasar durante su ejecución; aun así, sería un principio.)

El razonamiento lógico (correcto) y la planificación son subcampos importantes de la IA, y la IA clásica exploró, casi desde sus inicios, estrategias de razonamiento y planificación que usaban elementos de la lógica simbólica, como la deducción. Por ejemplo, un sistema de IA puede

implementar un silogismo y también planificar un algoritmo (reglas de la forma: $\{A, B, C, \dots\} \rightarrow G$, donde A, B y C son acciones que hay que realizar y G es el objetivo deseado). No se han producido grandes logros en el camino hacia la inteligencia artificial general usando esos métodos, pero incluso los científicos de IA modernos como Stuart Russell continúan insistiendo en que la lógica simbólica será un componente importante de cualquier sistema de inteligencia artificial general en el futuro —ya que la inteligencia trata, entre otras cosas, el razonamiento y la planificación.

Así, Aristóteles dio el pistoletazo de salida a los estudios formales sobre la inferencia hace miles de años. Hace unas pocas décadas, también ayudó a dar el pistoletazo de salida al trabajo en la IA. El razonamiento simbólico que usa reglas deductivas enlaza la inteligencia específicamente con el conocimiento, un prerrequisito del sentido común, que es lo que sigue faltando de manera casi completa en los sistemas de IA. John McCarthy, pionero de los inicios de la IA como fundador de la disciplina durante la conferencia de Dartmouth de 1956, se dio cuenta muy pronto de ello, y dedicó un esfuerzo continuado al desarrollo de los sistemas de conocimiento —sistemas que dependen, para razonar y actuar, de proposiciones sobre el mundo que se pueden representar por ordenador—. Todos los sistemas de conocimiento antiguos se encontraron con problemas abrumadores, pero instructivos. Quizá pueda volverse sobre algunos de esos problemas con la esperanza de obtener algún avance. Otros problemas, no obstante, parecen fundamentales. En concreto, son limitaciones inherentes al mismo razonamiento basado en reglas. La lógica deductiva resulta precisa porque nos proporciona certezas. Tal y como cabría esperar, la certeza le pone el listón bastante alto al mundo real, donde los sistemas de inteligencia artificial general (y la gente) deben demostrar su inteligencia.

LA DEDUCCIÓN: CÓMO NO EQUIVOCARSE NUNCA

Los lógicos (y los científicos informáticos) analizan las inferencias deductivas con sistemas de proposiciones que pueden ser verdaderas o falsas. Dicta la convención que todas las proposiciones que preceden a la

última dentro de una serie se llaman «premisas». La última proposición es una consecuencia de las premisas, y se conoce como «conclusión». Juntas, las premisas y la conclusión se denominan «argumento». Un buen argumento deductivo es una «apuesta segura», porque su conclusión es necesariamente verdadera. Aquí hay uno:

Cuando llueve, las calles se mojan.
De hecho, está lloviendo.
Por consiguiente, las calles están mojadas.

La conclusión es la inferencia que deberíamos extraer de las dos premisas. (En esencia, da respuesta a la pregunta: sin saber nada más, ¿qué se desprende de las premisas? La regla que se usa para inferir la conclusión es válida cuando la conclusión ha de ser verdadera siempre que las premisas lo sean también. La validez es un sello de «confianza» sobre esa regla, que preservará su veracidad siempre que nuestras premisas (o convicciones previas) sean verdaderas. Por consiguiente, el ejemplo de arriba es válido. Y utiliza una de las reglas deductivas más antiguas, que sigue recibiendo un nombre en latín: *modus ponens*. De forma casi simbólica:

Si P implica Q
Si P es verdad
 Q es verdad

Y, en formato completamente analizable (computable), tenemos su forma lógica:

$$\begin{array}{l} P \rightarrow Q \\ P \\ \hline Q \end{array}$$

Aquí, el conector « \rightarrow » tiene un significado específico, o semántico, que determina el valor en cuanto verdad de P y de Q . En la lógica deductiva, la regla se conoce como «condicional material», y nos garantiza que Q se deriva del carácter verdadero de P y de la regla $P \rightarrow Q$. (El rango de posibilidades verdaderas o falsas viene dado por una tabla de verdad, que mostraremos más tarde.)

Ahora pasa a considerar algunas modificaciones en nuestro argumento sobre la lluvia y las calles. En especial, ¿qué pasa si no llueve? En ese caso,

la regla no «se dispara». No se deduce nada. Pero la forma del argumento sigue resultando válida. Sigue siendo verdad que, cuando llueve, las calles se mojan. Cuando de hecho llueve, el argumento se vuelve «sólido» (y no tan solo válido). La solidez es verdad; una verdad real, en oposición a la verdad condicional de la validez. La solidez nos indica que las premisas son realmente «verdaderas». La solidez garantiza que los agentes inteligentes que usen la inferencia deductiva inferirán verdades a partir de verdades anteriores. La validez, por otro lado, solo garantiza que, crea lo que crea el agente inteligente, sus inferencias serán correctas formalmente (incluso si razona acerca de mentiras o falsedades). De hecho, los argumentos deductivos que son válidos, pero no sólidos, pueden introducir todo tipo de ridiculeces en el razonamiento deductivo. Por ejemplo:

Cuando llueve, los cerdos echan a volar.
Está lloviendo.
Por consiguiente, los cerdos están volando.

Se trata de un argumento estúpido, pero perfectamente válido, porque de nuevo utiliza el *modus ponens*, el modo de razonar a partir de una proposición hipotética. Por supuesto, la primera premisa es falsa. La segunda premisa también podría ser falsa si no está lloviendo en realidad. No obstante, por mucho que llueva, no podemos fiarnos de la primera premisa, porque no existe ninguna conexión entre la lluvia y que los cerdos vuelen —y, de todos modos, los cerdos no vuelan, con independencia del tiempo que haga o de cualquier otra cosa—. El argumento es válido, pero no sólido... y es completamente inútil.

He aquí una deducción sólida:

Todos los hombres son mortales.
Sócrates es un hombre.
Por consiguiente, Sócrates es mortal.

¿Cómo podría estar equivocada? No puede. La conclusión siguiente siempre presenta un 100 % de certeza. La deducción proporciona un patrón para el pensamiento «perfecto» y preciso de los seres humanos y las máquinas, y es en buena medida por ese motivo por lo que ha sido investigada de manera extensiva en matemáticas y en las ciencias, y se ha utilizado con éxito en numerosas aplicaciones de importancia en el ámbito de la IA. En sus inicios, por ejemplo, los sistemas de IA basados en la

deducción eran capaces de demostrar de manera automática teoremas reales (no «de juguete») en matemáticas. Un programa informático llamado Logic Theorist, creado por los pioneros de la IA Alan Newell, Herb Simon y Cliff Shaw, demostró teoremas lógicos de interés en una fecha tan temprana como 1956 sirviéndose de los *Principia Mathematica* de Bertrand Russell y Alfred North Whitehead, la obra fundacional de la lógica del siglo xx. Los sistemas de razonamiento automático que utilizan la deducción también se han aplicado al diseño de circuitos para placas base de ordenador, y a la tarea de verificar *software* y *hardware*, asegurándose de que el *software* no contiene errores ni contradicciones.¹ En tales casos, el enfoque deductivo resulta más sencillo y efectivo que los métodos modernos de la IA, que usan estadísticas y aprendizaje. Los primeros investigadores en IA sabían también que nuestro conocimiento se expresa a menudo de manera simbólica (como en el ejemplo de la lluvia), así que la deducción tiene sentido; es una elección obvia. Por desgracia, existen problemas bien conocidos a la hora de extender la inferencia deductiva a la inteligencia general.

PROBLEMAS DE CONOCIMIENTO

Con el paso de los años se han ido descubriendo numerosos problemas relacionados con la deducción. Quizá el más dañino haya sido que la deducción nunca añada conocimiento. Si sé que la gente es mortal (se muere) y que tal y Pascual es una persona, ya sé que tal y Pascual se va a morir. La deducción tan solo confirma la conclusión a la que una persona racional debería haber llegado a partir de las premisas proporcionadas, cosa que es fácil de ver en un silogismo simple porque el «conocimiento» se hallaba ya contenido en las proposiciones. La conclusión se limita a hacerlo explícito.

La deducción resulta extraordinariamente útil como defensa contra la posibilidad de que alguien infiera conclusiones delirantes o incorrectas a partir de un conjunto de proposiciones —por ejemplo, que insistan en que, según las premisas de la mortalidad humana y el hecho de que Sócrates fue un ser humano, deberíamos concluir que Alfa Centauri está hecha de queso. La deducción otorga a los agentes racionales una plantilla para que no «se

salgan de la senda», lo cual representa de manera evidente un buen primer paso para cualquier sistema de IA del que esperamos que acabe realizando inferencias inteligentes. Pero, usando solo la deducción, no llegaremos demasiado lejos. Por ejemplo, como respuesta a la teoría copernicana de que la Tierra gira alrededor del Sol y no al revés, los astrónomos de la vieja escuela ptolemaica podrían haber empleado un contraataque deductivo:

Si Dios creó el cielo, la Tierra se encontraría en el centro del cielo.
El cielo fue creado por Dios.
Por consiguiente, la Tierra se encuentra en el centro del cielo.

El argumento es válido pero, una vez más, esto solo nos indica que, si las premisas son de hecho verdaderas, la conclusión la sigue de manera necesaria. Todo el trabajo pesado recae en las preguntas empíricas acerca de la veracidad de las premisas. En la indagación sobre la mortalidad de Sócrates eso nos sale «gratis», por así decirlo, ya que en general todos coincidimos en que la gente se muere (por mucho que después vayan al cielo). Pero la generalización según la cual todo cielo creado por una divinidad tendría nuestro planeta en su centro parece tan debatible como cualquier otra afirmación bíblica o estética. Podríamos insistir en una interpretación alternativa de las Escrituras (es famosa la afirmación de Galileo de que Dios nos dice cómo funciona lo de ir al cielo, no cómo funciona el cielo). O podríamos, sobre todo si somos ateos o materialistas científicos, rechazar la veracidad de la segunda premisa sin pensárnoslo demasiado.

Por tanto, la deducción se vuelve inútil para la búsqueda de nuevos conocimientos; solo sirve para esclarecer convicciones enfrentadas cuando se han cometido auténticos errores en el razonamiento. Como todo el mundo sabe, es posible que los teóricos de la conspiración nunca cometan errores de razonamiento deductivo; es solo que adoptan como verdaderas unas premisas que para otras personas son dudosas o simplemente delirantes.

En otras palabras, todo sistema inteligente requerirá otros tipos de inferencia para centrarse en unas convicciones verdaderas (y útiles). No basta con la certeza deductiva de una conclusión inferida.

PROBLEMAS DE RELEVANCIA

La deducción presenta otras limitaciones que la vuelven inadecuada como estrategia para diseñar la inteligencia general. Una especialmente dañina involucra los factores de relevancia. La premisa «Si está lloviendo, los cerdos volarán» es falsa, porque los cerdos no vuelan, pero también es un ejemplo excepcionalmente malo de cómo decir algo relevante. La lluvia no tiene nada que ver con la cuestión de que los cerdos vuelen. Por otro lado, los aviones sí que vuelan, pero la premisa «Si está lloviendo, los aviones volarán» también resulta irrelevante. Podría ser verdadera (al menos en algunos casos), pero el hecho de que llueva no debería conducir a que tengamos alguna convicción acerca de que los aviones se eleven por los aires. Una vez más, la proposición ignora los factores de relevancia.

Parte del problema tiene que ver con la causalidad: la lluvia no hace que los aviones vuelen (aunque, en algunas circunstancias, puede hacer que se queden en tierra). Aquí depende de la manera en que queramos usar el conocimiento. La proposición «Si el termómetro llega al color rojo, ahí fuera hace calor» es verdadera. Pero, si queremos inferir una explicación posible para una ola de calor, el termómetro no nos ayuda en nada. La proposición es verdadera pero irrelevante. «Si el gallo canta es que está saliendo el sol» también es verdadera, pero si tuviéramos que preguntarle a un sistema de inteligencia artificial general por qué ha salido el sol y este nos saliera con el gallo, nos costaría atribuirle una gran inteligencia.

Piensa en este ejemplo, tomado del filósofo de la ciencia Wesley Salmon:

Los hombres que toman pastillas anticonceptivas con regularidad no se quedan embarazados.

Un hombre toma con regularidad las pastillas anticonceptivas de su esposa.

Por consiguiente, el hombre no se queda embarazado.²

De hecho, se trata de un argumento deductivo perfectamente sólido, que sigue el *modus ponens* con premisas verdaderas. Pero que el hombre no se quede embarazado no tiene nada que ver con las razones que se dan. Estas son irrelevantes, porque los hombres no se quedan embarazados de todos modos. El argumento no explica nada. Podemos imaginarnos un robot armado con una vasta base de datos y reglas razonando de esa manera, utilizando la deducción. En realidad, no hay nada *per se* erróneo, pero el

robot no ha comprendido nada... No sabe lo que es relevante y lo que es una estupidez.

Piensa en este ejemplo, más sutil:

Todo aquel que ingiera treinta gramos de arsénico morirá antes de veinticuatro horas.

Jonas ingirió treinta gramos de arsénico a la hora t .

Jonas murió antes de que pasaran veinticuatro horas desde la hora t .

Se trata de un argumento deductivo perfectamente correcto, pero no serviría para explicar la muerte de Jones si, por ejemplo, este hubiera ingerido el arsénico a la hora t y hubiera fallecido en un accidente de tráfico (quizá cuando acudía veloz al hospital) antes de expirar por envenenamiento. Aquí, de nuevo, el argumento se basa en una deducción correcta, pero resulta irrelevante. No nos cuenta nada. Es hasta engañoso. La relevancia, en otras palabras, a menudo presupone el conocimiento de la causalidad, donde un hecho produce en realidad un resultado, o hace que suceda algo.

Otro motivo por el que la deducción cae constantemente víctima de los problemas de relevancia es que, de manera invariable, hay muchas causas posibles para que ocurra algo en nuestra experiencia cotidiana (y en la ciencia). Accidentes como que se estrelle un avión, por ejemplo, se pueden analizar en general señalando una causa próxima (cercana) y una causa distante (alejada) que expliquen el desastre en conjunto. Por ejemplo, las recientes tragedias protagonizadas por la compañía Boeing. Después de que dos aviones Boeing 737 Max se estrellaran en un plazo de seis meses en 2018, los investigadores descubrieron un fallo en el *software* del sistema de estabilización, el Sistema de Aumento de Características de Maniobra (MCAS en sus siglas inglesas). Al rediseñarse el antiguo Boeing 737-800 se generó espacio para unos motores de mayor tamaño, pero a cambio de colocarlos algo más adelantados y ligeramente por encima de las alas. Eso condujo a que la velocidad vertical de ascenso durante el despegue fuera más elevada, lo cual podía provocar, bajo ciertas condiciones, problemas de sustentación. Los fallos de sustentación son malos —en potencia, catastróficos—, así que se dotó a los nuevos Max de un MCAS para que, cuando fuera necesario, este dirigiera el morro del avión hacia abajo y evitara el problema. Por desgracia, esa corrección del morro podía hacer que el Max se precipitara contra el suelo. Y eso hizo el MCAS: arrebató el

control del aparato a los pilotos en sendas tragedias que provocaron la muerte de 157 personas en Indonesia y de otras 189 en Etiopía.

La investigación consiguiente reveló fallas en el *software* que controlaba el MCAS, así que se identificó una causa próxima. Pero la investigación resultante de la anterior también recalcó el empeño con el que Boeing había puesto en servicio los Max a fin de competir en el ahorro de combustible que proporcionaban las aeronaves de Airbus, su gran rival —lo cual señaló una causa de fondo o distante—. También se descubrió que los pilotos del nuevo Max no habían recibido la formación adecuada. A ello contribuyó sin duda que, en el lanzamiento de mercadotecnia de la aeronave rediseñada, Boeing asegurara que el Max no requería que los pilotos ya familiarizados con el 737-800 tuvieran que someterse a un costoso proceso de aprendizaje. Así, aquellos trágicos accidentes pueden atribuirse a una multiplicidad de causas. Inferir el motivo por el que los Boeing 737 Max se estrellaron implica considerar diversas causas posibles, y quizá no haya una que por sí sola explique las catástrofes.

La deducción no puede dialogar con estos escenarios propios del mundo real. Al obligar a que las inferencias sean verdaderas más allá de cualquier duda, la deducción invariablemente pasa por alto aquello que puede ser verdad en contextos donde la relevancia viene determinada por una mezcla de factores que no son necesarios, pero que sí se mantienen operativos en ciertas situaciones. En el universo platónico de formas inmutables, los triángulos han de tener tres lados y algunas cosas son Verdad, con uve mayúscula. En la experiencia desordenada del mundo, observamos o analizamos pocas cosas parecidas al triángulo. Son más bien como el Boeing 737 Max —o igual que una conversación común y corriente (como ya veremos)—. La inteligencia —sea lo que sea— va más allá de las deducciones. Nosotros mismos somos sistemas cognitivos, y es evidente que no somos solo sistemas deductivos. Eso sugiere que, para tener éxito, la IA de nivel humano tampoco puede ser completamente deductiva.

Tras el fracaso de lo que los críticos catalogaron como « inteligencia artificial de la vieja escuela», que dominó la IA antes de la era contemporánea (hasta bien entrada la década de 1990), los científicos de IA abandonaron los enfoques deductivos para ponerse a inferir en masa. En efecto, a numerosos lectores jóvenes les parecerá extraño que los practicantes de la disciplina se tomaran alguna vez en serio cosas como las «reglas» y los enfoques deductivos. Así fue. Pero las limitaciones

devastadoras de la inferencia deductiva acabaron por condenar ese enfoque. Y, con la explosión de la red, las cantidades de datos disponibles para los llamados métodos superficiales o estadísticos hicieron que los sistemas deductivos o reglados parecieran menos útiles y más torpes. Un nuevo paradigma —un tipo de inferencia diferente— cobró prominencia en el trabajo serio sobre la IA. Se llama «inducción», y hablaremos de él a continuación.

EL PODER Y LOS LÍMITES DE LA INDUCCIÓN

La inducción implica la adquisición de conocimientos a partir de la experiencia. Por lo general, la experiencia es una interpretación en forma de observaciones —ver cosas—, aunque también puede provenir de cualquiera de nuestros cinco sentidos. (Tocar un fogón caliente es un ejemplo de inducción táctil.) A diferencia de la deducción, la forma general de la inducción pasa de las observaciones particulares a las hipótesis generales. La hipótesis inductiva cubre —esto es, explica— una observación. El mecanismo primario de la inducción es la enumeración: resulta difícil inducir los rasgos de una población de, pongamos, aves (por utilizar un ejemplo famoso) sin haber observado antes muchos ejemplos de aves. El carácter central de la enumeración desempeña un papel central en todas las versiones de la inducción, y será importante a la hora de comprender su naturaleza y sus limitaciones.

El poder de la inducción se debe no solo a que ayuda a organizar el mundo de las cosas en categorías a través de hipótesis (todos los objetos de X presentan la propiedad Y), sino que, además, otorga capacidades predictivas a los agentes que la utilizan.³ Si cada vez que se acaba un partido de la liga de béisbol las calles del centro se llenan a reventar de gente, podría inferir que eso volverá a pasar la próxima vez que se acabe un partido —lo cual es una predicción—. La inducción atrapa la idea cotidiana de que, al observar lo que acontece en el mundo, obtenemos la capacidad para explicarlo y predecirlo. Muchas de nuestras expectativas se basan en la inducción. Si alguien desplazara el picaporte de la puerta de entrada de tu casa diez centímetros a la izquierda, lo más probable es que fallaras al ir a cogerlo. Tienes una teoría implícita —esto es, una hipótesis—, basada en

los numerosos ejemplos previos en que lo has visto y cogido, y del lugar en el que se encuentra el picaporte.

La inducción tiene otras virtudes. Para comenzar es sintética, por tomar prestado el concepto de Kant; añade conocimiento. Puedo buscar en la red cuál es la hora punta del tráfico en la esquina de las calles Tercera y Mayor, pero si trabajo en el cruce de las calles Tercera y Mayor puedo mirar por la ventana. Esta última opción es una observación de primera mano que facilita mis inferencias inductivas y conforma expectativas y planes para cuando me vaya a casa. Por desgracia, la potente flexibilidad de la inducción (ligada a nuestros sentidos) también implica que no se trata de algo demostrable, ni de una verdad garantizada, como la deducción. El conocimiento cosechado por las observaciones es siempre provisional. ¿Por qué? Porque el mundo cambia. El futuro podría falsificar mis hipótesis inductivas. Es posible que mi coche se haya puesto en marcha mil veces sin problemas. Mañana por la mañana (cuando esté llegando tarde a una reunión: la ley de Murphy) quizá no sea así. Eso es la inducción. El cambio llega (o no, por desgracia), y la observación previa por sí sola no puede indicarnos cómo ni cuándo lo hará.

La fortaleza de la inducción, no obstante, radica en el hecho de que la inteligencia está ligada de manera importante al mundo que nos rodea. La ciencia moderna sería imposible sin su lealtad a la inducción como método de conocimiento a través de la experiencia.

Piensa de nuevo en las enumeraciones. En su forma más sencilla, la inducción requiere tan solo de la enumeración de las observaciones previas para llegar a la conclusión general o norma (ley). He aquí un argumento:

N cisnes observados han sido de color blanco [donde *N* es un número alto].
Por consiguiente, todos los cisnes son blancos.

O:

Toda la vida que hemos visto está basada en el carbono.
Por consiguiente, toda la vida está basada en el carbono.

Tal y como sugieren estos ejemplos, la enumeración simple (es decir, el conteo) de los rasgos o propiedades de algo a menudo conforma la base de nuestras pretensiones de conocimiento sobre el objeto como clase. Así, los cisnes son simplemente esas aves de color blanco; la vida no es más que el fenómeno que se deriva del carbono. En la ciencia (y en la vida) también

nos resulta útil contar una historia sobre el motivo por el que los cisnes podrían ser blancos, o por el que la vida podría tener una base de carbono, pero, en sentido estricto, las explicaciones con que se responde a los por qué quedan fuera del alcance de la inducción, la que enumera o la otra.

Sin embargo, es la simplicidad de la inducción lo que la vuelve tan útil como forma de inferencia. Cuanto más observo alguna propiedad en un objeto, más confianza tengo en que esa propiedad sea parte integral de ese objeto. Si no dejo de comprobar las pelotas de una bolsa y estas son siempre blancas, en algún momento comenzaré a confiar en una generalización del tipo «Todas las pelotas de esta bolsa son blancas». Pero, de nuevo, si no he comprobado hasta la última pelota, siempre existe la posibilidad de que mi inferencia inductiva sea errónea. La inducción resulta útil, pero no ofrece un conocimiento cierto.

He aquí otro tipo de generalización inductiva:

La proporción Q de una muestra de la población presenta la propiedad P .
Por consiguiente, la proporción Q de la población presenta la propiedad P .

La inducción entre una muestra y una población es bastante común en las investigaciones científicas, y a lo largo de los años se han desarrollado sofisticadas técnicas estadísticas para ayudar a que estas generalizaciones sean lo más sólidas y estén lo más libres de errores posible, dada la evidencia observable de la que se pueda disponer. En términos intuitivos, las generalizaciones inductivas también tienen sentido: si observo 75 bolas blancas y 25 bolas negras en una muestra, a falta de otra evidencia, debo esperar que haya 750 bolas blancas en una población de 1.000. La inferencia parece correcta, lo que pasa es que no existe una certeza sobre ella.

El muestreo aleatorio también se basa en una generalización debida a las observaciones. Pruébalo por ti mismo: lanza una moneda al aire varias veces y cuenta el número de caras y de cruces que te salen. Eso es un muestreo aleatorio (puesto que no puedes darle ningún sesgo al lanzamiento, la moneda es justa). Es posible que te salgan dos o tres caras seguidas. Aunque resulte muy improbable, puede ser incluso que te salgan cinco caras o cinco cruces seguidas. Pero, dado un muestreo lo bastante amplio, podrás generalizar diciendo que las posibilidades de que salga cara o cruz son mitad y mitad. De ahí que la generalización inductiva sea «La moneda lanzada al aire saldrá cara quinientas veces de cada mil», lo cual se

acerca bastante a la verdad. (La ley de los grandes números nos dice que, dado un muestreo lo bastante amplio, la probabilidad se aproximará a la probabilidad real: al cabo de un millón de lanzamientos, estaremos bastante cerca de ese 50-50 %.) He aquí otro ejemplo popular de generalización estadística a través de la inducción:

El 73 % de un muestreo aleatorio de votantes apoya al candidato *X*.
Por consiguiente, el candidato *X* obtendrá cerca del 73 % de los votos.

Existe la posibilidad de que el candidato *X* se vea envuelto en un escándalo antes de las elecciones, lo cual invalidaría esa inferencia inductiva. Pero, una vez más, en ausencia de un conocimiento mayor, podemos razonar de esta manera y extraer conclusiones sobre lo que esperamos que suceda.

La IA moderna se basa en análisis estadísticos, y por tanto depende de un marco inductivo, lo cual resulta de utilidad para numerosas aplicaciones comerciales. Por ejemplo, la IA puede ofrecer recomendaciones: un tipo de predicción basado en observaciones pasadas. He aquí otro ejemplo, que resultará familiar a cualquiera que tenga un buscador de contenido.

El 75 % de las noticias que lee el usuario *X* son crónicas políticas de carácter conservador en el sitio web *C*.
Por consiguiente, el usuario *X* querrá leer esta noticia en *C*.

Es posible que al usuario *X* también le guste de vez en cuando leer algún artículo del *New Republic*. Por desgracia, lo más probable es que el sistema que infiere las preferencias de *X* sirviéndose de la inducción lo ignore. Se trata de un inconveniente evidente a la hora de fiarse de las generalizaciones inductivas derivadas de la observación —son sustitutivos de un conocimiento más profundo (lo que es peor, tienden a esperar que el futuro tenga el mismo aspecto que el pasado).

David Hume, filósofo del siglo XVIII y el primer pensador que señaló los límites de la inducción, aportó a filósofos y científicos lo que ahora se conoce como el problema de la inducción. Tal y como él mismo lo expresó, confiar en la inducción requiere de nosotros la convicción de que «los casos de los que no tenemos experiencia se parezcan a aquellos de los que sí la tenemos». En otras palabras, la regla general inductiva que aplicamos tiene que ser ampliable a ejemplos que no hemos visto, y no existe ninguna garantía de que vaya a mantenerse. A diferencia de la deducción, en la

estructura de la inducción no hay nada que nos proporcione una certeza lógica. La inducción se limita a calcular que el mundo cuenta con ciertas características, y nosotros podemos examinarlo y sonsacarle el conocimiento (que creemos tener) sobre él.⁴

El problema de la inducción puede parecer una de esas preocupaciones menores a las que los filósofos se entregan por placer, pero en realidad los límites de la inferencia inductiva generan problemas constantes a los científicos en su búsqueda de teorías verdaderas. Hay ejemplos por doquier. Solíamos comernos solo la clara de los huevos porque la ciencia nutricionista nos había advertido contra los peligros de las grasas saturadas que se encontraban en la yema. Si pasamos algunas décadas a cámara rápida descubrirás que los científicos nutricionistas ahora nos animan a comer huevos con yema y todo, ya que ayudan a quemar grasas y levantan el ánimo; incluso protegen contra los problemas de corazón (el mismo motivo por el que nos provocaban ansiedad hace unas décadas). En un sentido muy real, podemos culpar a la inducción por esos cambios radicales tan vergonzosos. Estos suceden porque nuestras observaciones y comprobaciones nunca son completas. Las correlaciones pueden sugerir una causa subyacente en la que podamos confiar (un fragmento de conocimiento real), pero existe la posibilidad de que se nos haya pasado algo por alto mientras comprobábamos y observábamos las cosas que las afectaban. La correlación puede ser falsa o accidental. Podríamos haber estado buscando algo equivocado. La muestra podría haber sido demasiado pequeña o no representativa por razones que solo se vuelven aparentes más tarde. Es un problema habitual, y en el fondo se debe al espectro de la inducción y sus límites —al final resultó que los filósofos no nos estaban haciendo perder el tiempo.

En el fondo, toda inducción se basa en una enumeración. Quizá suene sospechosamente simple (o debería ser así): ¿es posible que, a fin de obtener teorías sobre el mundo, tengamos que limitarnos a contar ejemplos? En un sentido importante, sí. Una sola experiencia no permite la licencia de realizar inferencias inductivas. Si veo un orangután y conozco su aspecto, puedo clasificarlo. Pero, si aún no conozco qué animal es, tendré que observar muchos ejemplares antes de averiguar si el animal que he visto es un chimpancé peculiar o la cría de un Big Foot adulto. Tal y como dijo Hume, para inferir las causas necesitamos ver «correlaciones constantes», y para inferir categorías o tipos necesitamos contar con ejemplos enumerados.

(Como veremos, así es exactamente como funciona el aprendizaje automático.)

Por supuesto, el razonamiento inductivo se vuelve más complicado: en terrenos como la economía o las ciencias sociales, las inferencias estadísticas también son inductivas, pero para comprenderlas hay que saber muchísimo sobre teoría de probabilidades (y sobre economía y ciencias sociales). Y, en las ciencias, las nuevas inferencias inductivas se elevan de manera inevitable sobre otras más antiguas que los científicos han pasado a considerar sólidas y verdaderas. (Así que también tenemos que saberlo todo sobre esas teorías.) Pero, en el fondo, la inducción simplemente generaliza a partir de la observación de ejemplos. Cuando el origen de esas generalizaciones se puede explicar con una historia, con una causa o con un conjunto de causas, podemos confiar en que se haya adquirido un nuevo conocimiento..., aunque no sea cierto por necesidad, como con la deducción. Este cuenta con el apoyo de la observación y la comprobación.

La crítica que Hume realizó a la inducción fue sobre todo una crítica a la causalidad. La inducción no requiere que se conozcan las causas (de otro modo, no sería enumerativa). Si sabemos, por ejemplo, que el color del plumaje de las aves viene determinado en parte por las características de su hábitat, aunque todos los cisnes de Inglaterra sean blancos podemos esperar que se encuentren plumajes negros en cisnes de hábitats diferentes. Pero, a falta de una teoría, la inducción solo podrá indicarnos eso si nos ponemos a volar por todo el mundo y vamos observando a los cisnes en los diferentes lugares en los que viven. Las hipótesis que citan causas específicas son el objetivo de la observación, pero por desgracia los medios lógicos de la inducción no son adecuados para proporcionarlas. Hacen falta inferencias adicionales (y ahí la deducción puede ayudar, pero solo en parte).

La cuestión es la siguiente: la inducción entendida debidamente dentro del marco lógico de la inferencia es, aunque necesaria y habitual, bastante limitada. Y también suele ser malinterpretada, lo cual contribuye a un exceso de confianza general en el hecho de que la inducción garantice un conocimiento «científico» y de solidez empírica que nos libre de especulaciones rocambolescas. Sherlock Holmes, nuestro héroe detective, explica a veces su método como una serie de inducciones meticulosas, observaciones simples y claras en las que no se entromete ninguna opinión, ni idea, ni convicción. Asegura, ante un perplejo y fascinado Watson, que él se limita a «observar las cosas con detenimiento». Holmes conoce el valor

de la simple observación —cuanto más simple, mejor—, porque aquello que creemos saber puede impedir que veamos algo nuevo. Pero esa es solo una parte de la historia de la inteligencia. También tenemos que entender el sentido de lo que observamos. Holmes, igual que Dupin, soluciona crímenes recopilando observaciones de una manera novedosa. El diablo está en los detalles y los detalles se encuentran en la novedad, que no es ninguna inducción.

La inferencia inductiva nos sitúa inevitablemente delante de otro peligro, que el ojo crítico de Hume volvió memorable: los hechos recién descubiertos pueden sorprendernos. En entornos dinámicos como la vida cotidiana, la observación es abierta. Observaciones futuras pueden revelar lo que antes se encontraba oculto o nos resultaba desconocido: ¡sorpresa! Y nuestra propia confianza en la inferencia inductiva puede llevar a que nos cueste más evitar sus inevitables defectos y fracasos. Eso nos lleva a hablar de las fiestas navideñas o, al menos, del «pavo inductivo» de Bertrand Russell.

EL PAVO DE RUSSELL

Bertrand Russell fue uno de los filósofos e intelectuales públicos más famosos del siglo xx. Un lógico, matemático y activista social que en una ocasión pasó seis meses en la cárcel por protestar contra la entrada de Gran Bretaña en la primera guerra mundial. Más tarde, en los años cincuenta protestó contra la proliferación de armas nucleares. Sus intereses intelectuales también fueron una forma de protesta: le preocupaba que se pudiera usar el lenguaje para inventarse problemas y soluciones filosóficas, y pensó que el antídoto para esa filosofía ensoñada consistía en ligarla a los métodos de la ciencia.

Pero, tal y como el propio Russell señaló, la ciencia a menudo avanza sin unas reglas inferenciales claras. Para exponer la inferencia en la ciencia, por tanto, debemos revelar los errores de nuestro pensamiento acerca de la investigación científica y la búsqueda de la verdad en general. Ese fue el motivo por el que decidió poner la lupa sobre el problema de la inducción, que definió como uno de los «problemas filosóficos» fundamentales en un libro titulado precisamente *Los problemas de la filosofía*, y argumentó,

igual que sir Karl Popper, que la ciencia no acumula conocimiento coleccionando o enumerando hechos. En otras palabras, que no obtenemos conocimiento científico a través de la mera inducción. De hecho, la inducción en sí misma es irremediabilmente defectuosa.

Russell ofreció un ejemplo evidente y accesible: que observemos el sol elevarse cada mañana no nos ofrece ninguna prueba de que vaya a hacerlo de nuevo. Nuestra confianza en que el sol vaya a salir mañana no es más que un «hábito de asociación», tal y como lo denominó Hume. No es solo que la inducción resulte incompleta, sino que con toda seguridad no puede confirmar teorías científicas ni convicciones desde la enumeración de observaciones. Nuestra creencia de que sí lo hace genera todo tipo de distorsiones. La «falacia del jugador», por ejemplo, es una convicción tozuda que aparece en los ludópatas, según la cual la frecuencia pasada de un resultado nos comunica algo verdadero acerca de resultados futuros. La falacia puede apoyar la expectativa de que siga habiendo más de lo mismo o lo opuesto: que ha llegado el momento de que pase algo nuevo. Las rachas, cuando se juega a los dados, generan la ilusión de que de algún modo habrá una influencia sobre la siguiente tirada: la racha de buena suerte se mantendrá así (estoy que me salgo) o la racha de mala suerte tendrá que acabar (ya me toca). Cualquiera de los dos escenarios puede producirse, por supuesto, pero la moraleja importante es que la siguiente tirada de los dados es independiente de todas las tiradas anteriores. La buena racha continúa si los dados caen aleatoriamente de una manera y se rompe si caen de otra. Este es un ejemplo de nuestra ansiedad por aplicar la falacia inductiva incluso a los hechos más aleatorios.

Sin embargo, la mayor parte del mundo real no es aleatoria, y eso hace que cueste aún más arrancar de raíz la tendencia a ver patrones inductivos incorrectos —la verdad es que los patrones están «ahí fuera», pero no siempre podemos averiguar cuáles son los verdaderos solo a través de la observación—. Vemos regularidades y patrones por todas partes. Más allá de la ludopatía, este peculiar giro mental ayuda a explicar nuestra inclinación a generalizar a partir de la observación. Los cisnes son blancos. El sol volverá a salir. El ascensor siempre me está esperando en la planta baja a las 3:30 de la mañana. Por supuesto que hay generalizaciones en las que se puede confiar —las vemos por todas partes, y no es ilusorio hacerlo—, pero el problema de la inducción, tal y como señaló Russell, es que no tenemos ninguna base para inferir conocimientos basándonos solo en esas

generalizaciones. La ciencia debe confiar en estrategias inferenciales más profundas y poderosas. La inducción misma es fina como el papel.

Russell nos ofrece, a modo de ejemplo de los límites de la inducción, un ave de corral bien alimentada que resulta ser una fabulosa pensadora inductiva. He aquí una versión de su triste historia:

El pavo descubrió que, aquella primera mañana en la granja de pavos, le dieron de comer a las 9. No obstante, puesto que era un buen inductivo, no sacó conclusiones precipitadas. Esperó hasta haber reunido un amplio número de observaciones del hecho de que le daban de comer a las 9 de la mañana, y realizó esas observaciones bajo una amplia variedad de circunstancias, los miércoles y los jueves, los días de calor y los días de frío, los días de lluvia y los días secos. A diario, añadía otra proposición observacional a su lista. Al fin, su conciencia inductiva se dio por satisfecha y elaboró una inferencia inductiva para concluir: «Siempre me dan de comer a las 9 de la mañana». Pero, ay, esa conclusión se reveló falsa con una certeza absoluta cuando, el 24 de diciembre, en vez de darle de comer le cortaron el cuello. Una inferencia inductiva con premisas verdaderas ha conducido a una conclusión falsa.⁵

El pavo de Russell revela el disparate de crear «hábitos de asociación» sin tener un conocimiento más profundo de las regularidades que observamos. Pero el conocimiento a menudo es una convicción disfrazada: aquello que creemos saber puede ser erróneo.

Un segundo problema igual de dañino de la dependencia en la inferencia deductiva tiene que ver con la ausencia de conocimiento. Buena parte del mundo se halla oculta de manera misteriosa; está enfangada en lo aleatorio y lo caótico, o simplemente es demasiado compleja para que confiemos solo en la inducción. Me vienen a la cabeza los mercados financieros. Podemos intentar predecir el rendimiento de un paquete de acciones con todo tipo de técnicas sofisticadas, pero como cualquier agente de bolsa sabe bien, los resultados del pasado no son indicativos de los resultados del futuro. Y, si hemos de ser sinceros, buena parte de nuestra experiencia del mundo presenta esa cualidad frustrante. Sabemos que el ascensor se queda en la planta baja cuando no lo usa nadie y, por inducción, es posible que creamos que esté allí esperándonos si volvemos a casa pronto del trabajo, porque estaremos fuera del horario del resto de la gente. Pero alguien podría estar mudándose, o unos parientes de Minnesota han venido a visitar a fulano de tal, etc. Las reglas están hechas para romperse, y las expectativas, también.

Nuestras predicciones se ven constantemente frustradas porque el conocimiento que necesitamos para amplificar la inducción a menudo no existe o no se puede disponer de él. Puedo ver mil cisnes blancos en

Inglaterra y llegar a la conclusión de que «Todos los cisnes son blancos». Ese mismo año, durante un viaje por Australia, veo un cisne negro... y la inducción se va al garete. Buena parte de lo que creemos saber en realidad es provisional, está pendiente de exámenes ulteriores, y es esa sobredependencia de la inducción la que lleva a que los cambios parezcan sorprendentes. En las ciudades grandes del oeste de Estados Unidos, como Seattle, los conductores suelen reducir la velocidad o detenerse ante un semáforo en amarillo en vez de pisar a fondo el acelerador para dejarlo atrás. También tratan con deferencia a los peatones, en vez de sortearlos. Eso podría hacer que me sorprendiera el comportamiento de los conductores de Nueva York o de Bombay. Incluso bajo una misma normativa, la conducta difiere. Si me fío de los datos e inducciones de mis experiencias pasadas, es posible que me choquen por detrás o que me peguen un bocinazo.

Así que ¿por qué debería confiar únicamente en hechos pasados, con toda esta nueva información? ¿Qué sería lo más razonable e inteligente por mi parte? ¿Cómo debería tratar la nueva información? ¿Y si veo una tira de clavos atravesando la carretera, cosa a la que nunca me había enfrentado, o una fila de patos cruzándola, o señales de tráfico con las que no estoy familiarizado? Por desgracia, la respuesta en estos casos no es una mayor inducción, sino menos.

LA INDUCCIÓN FUNCIONA EN LOS JUEGOS, NO EN LA VIDA

El mundo real es un entorno dinámico, lo que quiere decir que se encuentra en cambio constante de maneras tanto predecibles como impredecibles, y que no podemos acotarlo con un sistema de reglas. Los juegos de mesa, no obstante, sí están acotados por sistemas de reglas, lo cual contribuye a explicar por qué los enfoques inductivos que aprenden a partir de la experiencia del juego funcionan tan bien. AlphaGo (o su sucesor, AlphaZero) utiliza un tipo de aprendizaje automático conocido como aprendizaje profundo para jugar al go, un juego de gran dificultad. Juega contra sí mismo, usando algo que se llama aprendizaje profundo por

refuerzo, e induce hipótesis acerca de los mejores movimientos que se pueden realizar sobre el tablero dada su posición y la del rival. El enfoque ha tenido un éxito fabuloso en «juegos de reglas conocidas, con dos participantes, discretos y observables», tal y como señala el científico de IA Stuart Russell.⁶ Es posible que Russell no haya estado pensando en el pavo del otro Russell, pero debería haberlo hecho: el problema real con los juegos que apuntalan la IA es que permiten la formación de hipótesis (generalizaciones nacidas de la experiencia) según unas reglas conocidas. Irónicamente, igual que la IA clásica antes, esas reglas no tienen aplicación en el mundo real, lo cual representa el quid de la búsqueda de una inteligencia general.

Los científicos informáticos que confían en los métodos inductivos suelen desdeñar y catalogar de irrelevantes los problemas de inducción señalados por Hume (o por Russell). Según la lógica, por supuesto que no existen garantías de corrección cuando se usa la inducción, pero sí que podemos «acercarnos bastante» a ella.

Esa respuesta yerra por completo el tiro. Un método catalogado como «probable, aproximadamente correcto» rige la formación de hipótesis para la IA estadística, como el aprendizaje automático, y se sabe que resulta efectivo para extirpar, con el tiempo, aquellas hipótesis que sean malas o falsas. Pero en realidad ese método es una extensión del argumento original de Hume, según el cual la inducción no ofrece ninguna garantía de corrección más que cuando se aplica a escenarios como el de los juegos, que disponen de reglas con las que encerrar las inferencias estadísticas. Una solución probable y aproximadamente correcta no altera el problema de la inducción en entornos dinámicos y ajenos al mundo de los juegos o a los laboratorios de investigación.

Los investigadores de IA son conscientes del problema de la inducción (explícita o implícitamente), pero este rara vez entra en las críticas al aprendizaje automático (o al aprendizaje profundo) porque, en esencia, estas se dedican a darle la vuelta al problema. Puesto que la inducción no funciona en entornos dinámicos, aceptan, la aplicamos en entornos controlables. Lo cual es como buscar las llaves debajo de una farola porque allí hay más luz que donde se han caído. Es cierto que los seres humanos han «solucionado» el problema de la inducción lo bastante bien como para utilizar la experiencia de manera efectiva en el mundo real (¿y dónde, si no?). Pero hemos solucionado ese problema no con una forma más potente

de inferencia inductiva, sino combinándola de algún modo con tipos más potentes de inferencia que contribuyan a la comprensión. El aprendizaje automático se da solo con inducción (tal y como comentaremos en el capítulo 11), así que los investigadores de la disciplina deberían mostrarse más escépticos de lo que son habitualmente sobre las perspectivas de la inteligencia artificial general.

REGULARIDAD Y FRAGILIDAD

La inducción permite a la inteligencia comportarse como un detector de regularidad. La IA estadística destaca a la hora de capturar regularidades a partir del análisis de datos, y ese es el motivo por el que las labores de reconocimiento visual de objetos, como la identificación de fotos de rostros humanos y mascotas, se cuentan entre sus éxitos. Los píxeles de un rostro se pueden distribuir y regular de manera que sean estudiados y clasificados. Sin embargo, puesto que esos sistemas aprenden a partir de las observaciones de unos patrones de impulso específicos, tienen problemas de fragilidad. Tal y como han señalado Gary Marcus, Ernest Davis y otros investigadores, incluso cambios en apariencia benignos como pasar el color de fondo del blanco al azul en las tareas de detección de objetos pueden llevar a que su rendimiento se degrade. Atestar las fotos con otras imágenes también tiene como resultado una degradación severa.⁷ La persona ignorará sin problemas que se añadan unas letras sin importancia en la zona roja de una señal de stop, pero cuando le presenten esta imagen alterada a un sistema de aprendizaje profundo, este la catalogará como una señal de límite de velocidad. Y hay ejemplos parecidos en el mundo real, incluyendo sistemas autónomos de navegación en coches sin conductor que han catalogado de manera errónea un autobús escolar como una máquina quitanieves, y un camión que giraba como un paso elevado.

El aprendizaje automático es inductivo porque adquiere conocimientos a partir de la observación de los datos. La técnica conocida como aprendizaje profundo es un tipo de aprendizaje automático —una red neuronal— que se ha revelado muy prometedor a la hora de reconocer objetos en fotos, de estimular el rendimiento de los vehículos autónomos y de participar en juegos en apariencia difíciles. Por ejemplo, el sistema DeepMind de Google

aprendió a jugar a diversos videojuegos clásicos de Atari y se celebró a bombo y platillo. Se proclamó que la suya era una inteligencia general, porque el mismo sistema era capaz de dominar diferentes juegos usando el llamado enfoque de aprendizaje por refuerzo profundo que propulsaba AlphaGo y AlphaZero. Pero la empresa emergente de IA Vicarious, por su parte, no tardó en señalar que cualquier cambio en apariencia inocuo sobre los juegos degradaba de manera manifiesta la actuación fabulosa del sistema. En Breakout, por ejemplo, el jugador mueve una pala de lado a lado sobre una línea base y va golpeando una pelota para lanzarla hacia arriba y contra un muro de varias capas de ladrillo. Con cada golpe destruye un ladrillo (y se acerca un poco más al «Breakout», a la evasión), pero cuando la bola rebota el jugador debe preocuparse de que no se le escape. Acercar la pala unos pocos píxeles a los ladrillos tuvo como resultado una degradación severa de su desempeño. «El sistema entero de DeepMind se derrumba», observaron Marcus y Davis en su crítica a la IA moderna. Y citaron una observación del pionero de la IA Yoshua Bengio según la cual las redes neuronales profundas «tienden a aprender regularidades estadísticas en el conjunto de datos en vez de conceptos abstractos de nivel superior».⁸

Algo que se suele ignorar o malinterpretar es que esos fracasos son fundamentales y no se pueden remendar con enfoques de aprendizaje más potentes y dependientes de la inferencia inductiva (basada en datos u observaciones). El problema aquí es el tipo de inferencia, no las especificidades de un algoritmo. Puesto que se necesitan muchos ejemplos para estimular el aprendizaje (en el caso del go, los juegos de ejemplo se elevan a varios millones), los sistemas son motores de inducción enumerativa glorificados, guiados por la formación de hipótesis dentro de los límites de las características del juego y de sus reglas. Esos mundos están cercados por reglas y son regulares —se trata de una especie de mundo en forma de curva de campana donde los mejores movimientos son los que con mayor frecuencia conducen a la victoria—. No se trata de ese mundo real que la inteligencia artificial general debe dominar, y que descansa más allá de los juegos de diseño humano y los centros de investigación. Esa diferencia lo es todo.

El pensamiento del mundo real depende de la detección sensitiva de la anormalidad o de las excepciones. Una calle urbana concurrida, por ejemplo, está llena de excepciones. Es uno de los motivos por los que no

hay robots paseándose por Manhattan (ni, por otro motivo relacionado con las excepciones, conversando con los seres humanos). Un robot manhattaniano no tardaría en tropezarse, provocaría un atasco de tráfico al aventurarse de manera desaconsejable a cruzar la calle, chocaría contra la gente o algo peor. Manhattan no es Atari ni el Go, y tampoco es una versión a escala de ellos—. Un «cerebro» de aprendizaje profundo sería (y es) un grave lastre en el mundo real, igual que cualquier sistema inductivo que pretenda reemplazar la inteligencia genuina. Si pudiéramos enseñar al pavo de Russell a que «jugara» a no convertirse en la cena, quizá aprendería a esfumarse el día antes de Navidad. Pero en ese caso no sería un buen pavo inductivista; dispondría de un conocimiento previo, suministrado por los seres humanos.

La IA estadística acaba así con un «problema de larga cola», donde los patrones habituales (en la cola gruesa de la curva de distribución) resultan sencillos, pero los raros (su cola larga) son difíciles. Por desgracia, algunas de las inferencias que realicen los sistemas de IA de inteligencia humana se encontrarán en la cola larga, no en el dulce punto de inducción de las regularidades localizables en los sistemas de mundo cerrado. De hecho, al centrarse en los éxitos «fáciles» que explotan las regularidades, la investigación en IA corre el peligro de alejarse de manera colectiva del avance hacia una inteligencia general. Ni siquiera estamos realizando un progreso gradual, porque, en la práctica, trabajar los problemas sencillos implica descuidar los reales (las llaves no están cerca de las farolas). Por sí mismas, las estrategias inductivas generan una esperanza falsa.

Que una foto quede mal catalogada en Facebook o que Netflix nos recomiende una película aburrida quizá no represente un problema demasiado grande con la inducción dirigida por datos, pero los coches sin conductor y otras tecnologías de carácter más crítico sin duda pueden meternos en un lío. Oren Etzioni, responsable del Instituto Allen de Inteligencia Artificial, afirma que el aprendizaje automático y los macrodatos son «modelos estadísticos de alta capacidad».⁹ Eso habla de unas ciencias informáticas impresionantes, pero no de una inteligencia general. Las mentes inteligentes suman la comprensión a los datos, y pueden encontrar el sentido que lleve a apreciar los puntos de falla y las anomalías. Los datos y el análisis de datos no son suficientes.

EL PROBLEMA DE LA INFERENCIA EN TÉRMINOS DE CONFIANZA

En su esclarecedora crítica a la inducción que se utiliza en las predicciones financieras, el ex agente de bolsa Nassim Nicholas Taleb divide los problemas de predicción estadística en cuatro cuadrantes, siendo sus variantes: primero, si la decisión que se ha de tomar es simple (binaria) o compleja, y, segundo, si la aleatoriedad del asunto es «mediocre» o extrema. Los problemas del primer cuadrante reclaman decisiones simples en relación con una distribución de probabilidad de cola fina. Los resultados son relativamente fáciles de predecir en términos estadísticos, y los sucesos anómalos, cuando se dan, tienen un impacto pequeño. Los problemas del segundo cuadrante son fáciles de predecir, pero, cuando sucede algo inesperado, las consecuencias son importantes. Los problemas del tercer cuadrante implican decisiones complejas, pero consecuencias manejables. Y luego están los problemas «pavo», en el cuarto cuadrante, que implican decisiones complejas emparejadas con distribuciones de probabilidad de cola gruesa y consecuencias de alto impacto. Piensa en los cracs del mercado de valores. Taleb apunta al exceso de confianza en la inducción como uno de los factores clave que han exacerbado el impacto de esos sucesos. No se trata solo de que nuestros métodos inductivos no funcionen, sino de que, al fiarnos de ellos, no logramos utilizar enfoques mejores, y eso tiene consecuencias potencialmente catastróficas. En efecto, nos obsesionamos con el pensamiento maquinal, cuando el análisis del pasado no nos ayuda en nada. Este es uno de los motivos por los que la superinteligencia inductiva llevará a resultados estúpidos. Tal y como dice Taleb en broma, es importante saber cómo «no convertirse en un pavo».¹⁰

Por supuesto, la predicción cuenta con otros límites que no se pueden resumir con pulcritud desnudando los puntos ciegos de la inducción. Los cisnes negros, al fin y al cabo, son raros, igual que los cracs del mercado de valores y las guerras (y las innovaciones) de importancia. Se nos puede perdonar que usemos la inducción para ayudar a iluminar unas posibilidades que de todos modos son opacas y muy impredecibles, pero no que intentemos reemplazar nuestro entendimiento exclusivamente con datos y estadísticas. En algunos casos, como en los sistemas naturales caóticos (pongamos, los sistemas con turbulencia), ahora sabemos que habrá

limitaciones inherentes a la predictibilidad cuando se use cualquier tipo conocido de método inferencial. La inducción puede sugerir que el pasado será parecido al futuro, pero la teoría del caos nos dice que no será así —o, al menos, que no podemos determinar en qué sentido se parecerán—. En algunos casos, aunque incompleto, el análisis estadístico es todo lo que tenemos.

CAUSA PROBABLE

El ganador del premio Turing Judea Pearl, un célebre científico informático que ha dedicado su vida al desarrollo de métodos computacionales que se muestren efectivos con el razonamiento causal, sostiene en su obra de 2018 *El libro del porqué* que el aprendizaje automático nunca podrá ofrecernos una comprensión real porque el análisis de los datos no abarca el conocimiento de la estructura causal del mundo real, esencial para la inteligencia. La «escalera de causalidad», tal y como la llama él, parte de la asociación de puntos de datos (ver y observar), y continúa con una intervención en el mundo (hacer), lo cual requiere el conocimiento de las causas. A continuación, sigue avanzando hacia el pensamiento contrafactual, como la imaginación, el entendimiento y las preguntas de tipo: ¿y si hubiera hecho algo de manera diferente?

Los sistemas de IA que usan métodos de aprendizaje automático —y muchos animales— se encuentran en el peldaño más bajo, el de la asociación. En ese primer nivel, buscamos regularidades en nuestras observaciones. Es lo que hace la lechuza cuando observa los movimientos de la rata y averigua el lugar donde es probable que esté el roedor un momento después, y es lo que hace un programa de ordenador para jugar al go cuando estudia una base de datos de millones de partidas para averiguar qué movimientos están asociados a un porcentaje mayor de victorias.¹¹

Aquí, Pearl nos hace el favor de relacionar observaciones y datos.¹² También señala que el ascenso por esa escalera implica diferentes tipos de pensamiento (o, de manera más específica, de inferencia). La asociación no guarda proporción con el pensamiento causal, ni con las imaginaciones. Podemos reestructurar el problema del salto entre la inteligencia artificial y la inteligencia artificial general como, precisamente, el problema de

descubrir nuevas teorías que permitan subir por la escalera de Pearl (o, en el marco contemporáneo, el problema de pasar de la inducción a otros tipos de inferencia más potentes).¹³

UN MANUAL BÁSICO DE SENTIDO COMÚN

Es posible que tus padres, o tu pareja o un amigo, te hayan acusado alguna vez de no tener sentido común, pero ánimo: tienes mucho más que cualquier sistema de IA, pero de lejos. Como bien sabía Turing, el sentido común es lo que permite que dos personas mantengan una conversación cotidiana. El problema del sentido común y de la comprensión de un lenguaje particular, que lo requiere, ha sido un motivo de inquietud notable entre los investigadores de IA desde la aparición de la disciplina. Y al final se está poniendo de manifiesto que el bombo y platillo que han acompañado el aprendizaje automático no nos están acercando nada a él. Los investigadores están reconociendo cada vez más este hecho, y no podría haber llegado en un mejor momento. Marcus y Davis se preguntan, si los ordenadores son tan listos, ¿por qué no pueden leer ni conducirnos por «El sentido común y el camino hacia el entendimiento profundo»? (uno de los capítulos de su libro).¹⁴ Stuart Russell encabeza su listado de «Avances conceptuales que están por llegar» con los aún misteriosos «lenguaje y sentido común».¹⁵ Pearl también entiende que la comprensión del lenguaje sigue sin resolverse (y ofrece su propio «minitest de Turing», que requiere el entendimiento de la causalidad).¹⁶

Así que, para progresar en la IA, tenemos que ir más allá de la inducción. (Si en una escalera metafórica te encuentras en el peldaño de la asociación, mira hacia arriba.) Hagámoslo a continuación —o al menos comencemos a hacerlo—. En el camino hacia la necesidad de una inferencia abductiva, primero deberíamos ir a lo específico; en particular, al aprendizaje automático y su fuente de entrada, el *big data*.

Capítulo 11

El aprendizaje automático y el *big data*

Aprender consiste en «mejorar unas prestaciones basándose en la experiencia».¹ El aprendizaje automático consiste en lograr que los ordenadores mejoren sus prestaciones basándose en la experiencia.

Esta definición del subcampo de la IA conocido como aprendizaje automático goza de una amplia aceptación y no resulta especialmente controvertida. Se ha mantenido esencialmente inalterada desde los trabajos tempranos sobre algoritmos de aprendizaje en los albores de la disciplina. Tom Mitchell, científico informático de Carnegie Mellon e investigador sobre el aprendizaje automático desde hace mucho, ofreció una definición ligeramente más detallada en su obra *Machine Learning* [«Aprendizaje automático»], de 1997: «Se dice que un programa informático ha aprendido de la experiencia E con respecto a algún tipo de tareas T y una medición del rendimiento P cuando su rendimiento con las tareas de T , medido por P , mejora con la experiencia E ».² En otras palabras, el aprendizaje automático es el tratamiento informático de la inducción: la adquisición de conocimiento a partir de la experiencia. El aprendizaje automático no es más que una inducción automatizada, así que no debería sorprendernos que los problemas con la inferencia inductiva impliquen otros problemas para el aprendizaje automático. Desarrollar esos problemas inevitables es el objetivo de este capítulo.

Hay dos tipos principales de aprendizaje. Cuando los seres humanos etiquetan el dato de entrada para señalar el resultado deseado, se llama «aprendizaje supervisado». Por el contrario, cuando el sistema analiza los patrones que pueda haber en los datos tal y como son, se llama «aprendizaje no supervisado». También hay un término medio. El «aprendizaje semi-supervisado» se inicia con una semilla, o pequeña parte de datos, que ha sido preparada por los seres humanos, y a continuación la va proyectando cada vez sobre una mayor cantidad de datos sin supervisión.

En los últimos años, los científicos de IA se han centrado ampliamente en un tipo de aprendizaje automático específico, el llamado «aprendizaje profundo», que ha ofrecido resultados impresionantes como enfoque de aprendizaje supervisado. A continuación comentaré el aprendizaje supervisado con cierto grado de detalle, además del aprendizaje profundo y sus aplicaciones. Como puedes suponer, el aprendizaje supervisado es un enorme manto que aglutina diferentes tipos de aprendizaje; los exploraré para ofrecer una visión general de los problemas que encuentra la IA.

Un tipo común de aprendizaje supervisado es la clasificación, que ha sido ampliamente tratada en centros de investigación y aplicaciones comerciales. Por ejemplo, los clasificadores aprendidos filtran el correo basura. El resultado de salida es un sí o no binario: el mensaje o bien es correo basura o no lo es. Por lo general, el sistema clasificatorio del correo basura se encuentra supervisado por el usuario de la cuenta de correo, que marca los mensajes entrantes como correo basura y los manda a la carpeta de correo no deseado o a la basura. En un segundo plano, el sistema de aprendizaje automático etiqueta aquellos mensajes que sean ejemplos positivos de correo basura. Cuando ya ha reunido una cantidad suficiente de ejemplos, el sistema aprende por sí mismo usando esos y otros mensajes entrantes, y crea un circuito de retroalimentación que converge en la diferencia entre los mensajes aceptables y los que son correo basura.

El filtro del correo basura es uno de los primeros ejemplos de la utilidad del aprendizaje automático en la red. Los algoritmos bayesianos ingenuos y otros clasificadores sencillos de probabilidad asignan puntuaciones numéricas a las palabras del mensaje indicando si son correo basura o no, y es el usuario quien facilita las categorías de correo deseado y no deseado. Al final, el clasificador acaba teniendo una hipótesis o modelo de correo basura basado tan solo en el análisis de las palabras de los mensajes. Los mensajes futuros se filtrarán de manera automática y los mensajes no

deseados irán a parar a la carpeta de correo basura. Hoy en día, los clasificadores de correo basura utilizan muchísimo conocimiento suministrado por el hombre —pistas sobre lo que representa un correo no deseado, como ciertas palabras en el lema, términos y frases «basuriles» conocidos, etc—. Los sistemas no son perfectos, en buena medida por culpa del constante juego del gato y el ratón que se da entre los proveedores de servicios y los creadores de correo basura, que no dejan de probar enfoques nuevos y diferentes para engañar a los filtros entrenados.³

La detección del correo basura no constituye un ejemplo demasiado sexi de aprendizaje supervisado. Los sistemas de aprendizaje profundo modernos también recurren a la clasificación para tareas como el reconocimiento de imágenes o de objetos visuales. Los populares concursos de ImageNet presentan ante sus participantes una tarea de aprendizaje supervisado a gran escala, con millones de imágenes que ImageNet ha descargado de sitios web como Flickr para usarlas como entrenamiento y poner a prueba la precisión de los sistemas de aprendizaje profundo. Todas esas imágenes han sido etiquetadas por personas (que han ofrecido sus servicios al proyecto a través de la interfaz Mechanical Turk de Amazon) y los términos que aplicaron han servido para establecer una base de datos estructurada de palabras inglesas conocida como WorldNet. Cada subconjunto seleccionado de las palabras de WorldNet representa una categoría que hay que aprender, usando nombres comunes (como «perro», «calabaza», «piano», «casa») y una selección de objetos más desconocidos (como «Scottish terrier», «mono rojo», «flamenco»). El concurso consiste en ver cuál de los clasificadores de aprendizaje profundo que compiten es capaz de etiquetar más imágenes de manera correcta, tal y como lo hicieron las personas. En los concursos de ImageNet se utilizan más de mil categorías, así que la tarea supera con creces el problema de sí/no que se presenta a los detectores de correo basura (o en cualquier otra labor de clasificación binaria, como la que se limita a etiquetar si la imagen pertenece a una cara humana o no). Participar en esta competición implica realizar una tarea masiva de clasificación usando los datos de píxel como estímulo.⁴

En las aplicaciones que procesan un lenguaje natural se suele usar la clasificación secuencial. Se trata las palabras como si tuvieran un orden definido, una secuencia. La clasificación documental o textual puede usar un enfoque más simple, sin ordenación —como un modelo «bolsa de

palabras», BOW según su acrónimo inglés—, pero esa información adicional donde las palabras se ven como un texto ordenado suele mejorar el rendimiento de la clasificación de textos. Por ejemplo, las palabras que aparecen en el título y en el primer párrafo nos proporcionan a menudo pistas bastante potentes sobre el significado o el tema del artículo. La clasificación de textos puede explotar esos rasgos a la hora de autoetiquetar artículos con categorías como «CIENCIA», «NEGOCIOS», «POLÍTICA» y «DEPORTE». La clasificación de textos es otro ejemplo de aprendizaje supervisado, porque se inicia cuando los seres humanos etiquetan artículos debidamente con sus categorías y proporcionan un estímulo de entrada al sistema de aprendizaje. Igual que la colección de fotos correctamente etiquetadas de ImageNet, también hay corpus amplios, o conjuntos de datos creados por personas al comentar colecciones de textos, cuyos metadatos sobre temas y demás rasgos resultan útiles para la formación de los sistemas de aprendizaje supervisado en labores de procesamiento del lenguaje.

El aprendizaje automático supervisado se encuentra tras una buena parte de la red contemporánea. Por ejemplo, posibilita la personalización de noticias y otras fuentes de contenido. Cuando el usuario hace clic principalmente sobre noticias de política, un algoritmo de aprendizaje supervisado que se está ejecutando en segundo plano (pongamos que en los servidores de Facebook) le presentará cada vez más noticias sobre política. Otros enfoques más sofisticados clasifican esas noticias políticas según su punto de vista; ofrecen noticias de corte más conservador o liberal al usuario cuya tendencia haya sido identificada por el sistema, e incluso las clasifican según sus sentimientos —como cuando un sistema clasifica los textos de opinión como positivos o negativos, igual que las críticas de cine.

Junto con la clasificación, los enfoques de aprendizaje supervisado se usan para etiquetar de manera automática los objetos individuales de una secuencia en vez de la secuencia entera, como sucede con la clasificación de imágenes y de textos. Es lo que se conoce como aprendizaje secuencial. Un ejemplo simple (aunque aburrido) de este es el etiquetado gramatical, donde una secuencia de palabras como «la vaca marrón» se etiqueta según sus funciones: «La / AD vaca / NC marrón / ADJ», refiriéndose esas etiquetas a sus caracteres respectivos de artículo determinado, nombre común y adjetivo. El aprendizaje secuencial no se sirve de las reglas de la lingüística para facilitar a los programas un conocimiento sobre las partes gramaticales del discurso; en cambio, las personas se limitan a etiquetar las

palabras de las frases con su categoría gramatical correcta y proporcionan esos datos de preparación humana como entrada del algoritmo de aprendizaje. Las máquinas resolvieron hace mucho tiempo el problema del etiquetado gramatical; el suministro de decenas de miles de frases marcadas alcanza un rendimiento de nivel humano sobre los datos sin examinar, es decir, cualquier frase que no se haya utilizado durante el entrenamiento. Otro problema muy explorado dentro del procesamiento del lenguaje es el del reconocimiento de entidades, donde el sistema de aprendizaje supervisado predice las entidades que aparecen en un texto, como menciones a personas, lugares, momentos, empresas y productos. La frase «El señor Smith informó de que XYZ Co. vendió más de diez mil aparatos durante el primer trimestre» podría etiquetarse como «El señor Smith / PERSONA informó de que XYZ Co. / COMPAÑÍA vendió más de diez mil / NÚMERO aparatos / PRODUCTO durante el primer trimestre / FECHA».

La clasificación secuencial también se puede usar para realizar predicciones de series temporales, donde los objetos previos permiten predecir cuál será el objeto siguiente. Los sistemas de reconocimiento de voz como Siri son un tipo de predicción de series temporales, igual que los populares sistemas de voz a texto. La predicción de series temporales tiene importantes aplicaciones en tareas tan complejas como los diagnósticos médicos, la planificación industrial y el precio de las acciones, entre otros.

El aprendizaje supervisado es la causa de casi todos los éxitos de importancia que ha cosechado el aprendizaje automático hasta la fecha, incluyendo el reconocimiento de imágenes o de voz, la navegación autónoma de vehículos sin conductor y la clasificación de textos y estrategias de personalización en la red. El aprendizaje no supervisado presenta la virtud de requerir una preparación mucho menor de los datos, ya que no son los seres humanos quienes añaden las etiquetas a los datos de entrenamiento. Pero una consecuencia directa de esa pérdida de la «señal» humana es que los sistemas no supervisados llevan un gran retraso en comparación con sus primos supervisados a la hora de realizar las tareas del mundo real. El aprendizaje no supervisado resulta útil en tareas abiertas como la de permitir que el ser humano visualice los grandes conglomerados de datos que generan los algoritmos de aprendizaje no supervisado. Pero, ya que la mayor parte del bombo relacionado con el aprendizaje automático —y en especial con el aprendizaje profundo— está relacionada con el aprendizaje supervisado, voy a centrar la discusión principalmente en

él. Sin embargo, no olvidemos que todas las limitaciones de origen inductivo que presentan los enfoques de aprendizaje supervisado aparecen con mayor fuerza incluso en el aprendizaje no supervisado. Al centrarnos en el aprendizaje supervisado, estamos prestando atención al mejor y más potente de los casos.

EL APRENDIZAJE AUTOMÁTICO COMO SIMULACIÓN

Desde un punto de vista conceptual y matemático, el aprendizaje automático es intrínsecamente una simulación. Los diseñadores de cada sistema de aprendizaje automático examinan un problema de uso intensivo de datos, y si existe algún tratamiento posible de aprendizaje automático lo consideran «bien definido». Asumen que alguna de sus funciones puede simular una conducta del mundo real o un sistema real. Se considera que el sistema real cuenta con un patrón oculto que da pie al resultado observable en los datos. La cuestión consiste en no recopilar de manera directa el patrón oculto —lo cual obligaría a comprender algo más que los datos—, sino en simular ese patrón oculto analizando las «huellas» que haya dejado en los datos. Es una distinción importante.

Piensa en otra tarea del procesamiento de lenguaje, la que se conoce como «etiquetado de roles semánticos». En ella, los diseñadores del algoritmo de aprendizaje desmontan el significado de las frases en términos de preguntas comunes del tipo quién, para quién, qué y cuándo. La función del algoritmo de aprendizaje consiste en tomar frases de ejemplo como entrada y generar el resultado de un conjunto de etiquetas que respondan a esas preguntas identificando los roles semánticos expresados en la frase. Por ejemplo, la frase puede contener a un agente que realiza una acción, un tema (el objeto involucrado en la acción) y un beneficiario de esta, y etiquetarse así: «John / AGENTE tiró / ACCIÓN la pelota / TEMA a Lizzy / BENEFICIARIO». En todos esos casos, el enfoque del aprendizaje automático implica asumir un comportamiento que se da pero que es desconocido, y utilizar un enfoque de aprendizaje para imitarlo de la mejor forma posible, descubriendo una función f . Como resultado del

entrenamiento al que se somete al sistema se genera una f como modelo o teoría del comportamiento que se da en los datos. Este modelo puede registrar roles semánticos, o entidades, o partes gramaticales, o imágenes de peces dorados; todo depende de la tarea de aprendizaje—. El aprendizaje automático es, por naturaleza, la simulación de un proceso que resulta demasiado complejo o que es desconocido, en el sentido de que no se dispone de unas reglas fáciles de programación, o de que entenderlo correctamente requeriría de un esfuerzo demasiado grande. A veces, un aprendizaje no supervisado revela la existencia de patrones en los datos. Pero son los seres humanos quienes identifican ese patrón después de un análisis; el algoritmo no sabe buscarlo. Y, si supiera, estaríamos hablando ya de un aprendizaje supervisado.

La mayoría de nosotros conocemos las funciones por las clases de matemáticas en el colegio, y su ejemplo más clásico es de tipo aritmético: $2 + 2 = 4$ es una ecuación cuyo operador, el símbolo de la suma, es técnicamente una función. Las funciones generan respuestas únicas según su entrada: así, la función de suma nos devuelve un 4 para $2 + 2$ (y nunca un 5, salvo en las novelas de George Orwell). Los primeros científicos de IA asumieron que se podrían resolver numerosos problemas del mundo real suministrándoles reglas que equivalieran a funciones de salida conocida, como con las sumas. Sin embargo, resultó que la mayoría de los problemas de interés para los investigadores en IA tienen funciones desconocidas (si es que existe alguna función relacionada con ellos). Por ello, ahora tenemos el aprendizaje automático, que busca aproximarse o simular esas funciones desconocidas. Este carácter «falso» del aprendizaje automático pasa desapercibido cuando el rendimiento del sistema se acerca de manera notable al de los seres humanos, o lo mejora. Pero la naturaleza imitativa del aprendizaje automático se ve expuesta con rapidez cuando el mundo real se separa de la simulación aprendida.

Este hecho tiene una importancia enorme, y se pasa por alto demasiado a menudo en los debates sobre el aprendizaje automático. He aquí otro hecho: los límites del mundo de un sistema de aprendizaje automático quedan precisamente establecidos por los conjuntos de datos que se le proporcionen durante su entrenamiento. El mundo real no deja de generar conjuntos de datos: veinticuatro horas al día, siete días a la semana, a perpetuidad. Por ello, cualquier conjunto de datos dado es solo una fracción muy pequeña de tiempo que representa, en el mejor de los casos, una evidencia parcial del

comportamiento de los sistemas del mundo real. Ese es uno de los motivos por los que la larga cola de acontecimientos improbables resulta tan problemática: el sistema no cuenta con una comprensión verdadera del sistema real (en comparación con el simulado). Esto es de una importancia tremenda para los debates sobre el aprendizaje profundo y la inteligencia artificial general, y plantea una serie de consideraciones problemáticas sobre cómo, cuándo y hasta qué punto deberíamos confiar en unos sistemas que técnicamente no comprenden los fenómenos que están analizando (salvo por lo expresado en sus conjuntos de datos durante el entrenamiento). Volveremos sobre estos temas en capítulos posteriores, ya que son capitales para comprender el paisaje del mito.

Existen al menos dos problemas con el aprendizaje automático como camino en potencia hacia la inteligencia general. Uno, que ya hemos tocado, es que se puede aprender de manera exitosa, al menos durante un tiempo, sin la menor comprensión del tema. Un sistema entrenado puede predecir resultados, entendiendo en apariencia el problema, hasta que un cambio o suceso inesperados hacen que la simulación se vuelva inútil. De hecho, las simulaciones que fracasan, cosa que sucede muy a menudo, pueden volverse peor que inútiles: pensemos en el uso del aprendizaje automático a la hora de conducir, y que contar con sus predicciones automatizadas provoque una falsa confianza. Es algo que pasa por todas partes; el mundo real, tan desordenado, siempre está alterando su rumbo. Se cambia de tema de conversación. Las acciones siguen una tendencia al alza, pero acto seguido un suceso exógeno, como una reestructuración corporativa, un terremoto o una inestabilidad geopolítica, hace que se desplomen. Es posible que Joe se pirre por los blogueros conservadores hasta el día en que su amigo Lewis le sugiera una revista digital inclinada hacia la izquierda, que su fuente de noticias había descartado y le había ocultado por completo. Es posible que Mary adore los caballos hasta que Sally, el suyo, se muere y ella desarrolla una pasión por el zen. Etcétera. La verdad es que el nombre de «aprendizaje automático» resulta poco apropiado, ya que los sistemas no aprenden en el sentido en que lo hacemos nosotros, adquiriendo una valoración del sentido del mundo cada vez más sólida y profunda. Son más bien como curvas de campana del aprendizaje: simulaciones basadas únicamente en datos de aquello que experimentamos de una manera directa en el mundo real.

El sentido común sirve de mucho a la hora de comprender las limitaciones del aprendizaje automático: nos dice que la vida resulta impredecible. Por ello, la única crítica dañina de verdad que se le puede hacer al aprendizaje automático es que mire hacia atrás. Al basarse en observaciones procedentes de conjuntos de datos —esto es, observaciones previas— puede revelar patrones y tendencias que nos sean de utilidad. Pero todo aprendizaje automático es una fracción de tiempo procedente del pasado; ante un futuro abierto en el que los cambios son deseables, los sistemas han de pasar por un nuevo entrenamiento. El aprendizaje automático solo puede ir a la zaga del flujo de nuestra experiencia, simulando (lo que esperamos que sean) regularidades útiles. Es la mente —no la máquina— la que marca el camino.

EL APRENDIZAJE AUTOMÁTICO COMO UNA IA DÉBIL

La naturaleza imitativa del aprendizaje automático también ayuda a explicar por qué se encuentra permanentemente atascado en aplicaciones definidas de forma débil, y por qué sus progresos hacia una inteligencia artificial general son pequeños o inexistentes. Los problemas bien definidos que existen en el procesamiento de lenguaje natural, como la clasificación de textos, el etiquetado gramatical, el analizador sintáctico y el reconocimiento de correo basura, entre muchos otros, deben ser analizados de manera individual. Esos sistemas tienen que ser objeto de un rediseño amplio, han de ser transferidos para que puedan solucionar otros problemas, aun cuando sean similares. Resulta irónico que a esos programas los llamemos «aprendices» porque, para nosotros, el significado del verbo «aprender» implica en esencia huir de los recursos limitados para alcanzar una comprensión más general de las cosas del mundo. Pero los sistemas que juegan al ajedrez no juegan al go, de mayor complejidad. Ni siquiera los sistemas de go juegan al ajedrez. Incluso el muy publicitado sistema Atari del DeepMind de Google generaliza solo entre diferentes juegos de Atari, y ni siquiera logró aprender a jugar a todos ellos. Los únicos que se le dieron bien fueron aquellos que seguían unos parámetros estrictos. Los sistemas de

aprendizaje más potentes son mucho más frágiles y limitantes de lo que podríamos suponer. Pero tiene sentido, porque los sistemas no son más que simulaciones. ¿Qué otra cosa podíamos esperar?

Los problemas con la inducción que hemos comentado más arriba no se derivan de la experiencia *per se*, sino del esfuerzo por fundamentar el conocimiento y la inferencia solo en esa experiencia, que es precisamente lo que hacen los métodos de aprendizaje automático para llegar a la IA. No debería sorprendernos, por tanto, que esos métodos —y los que están centrados en datos— padezcan también todos los problemas de la inducción. Los datos solo son hechos que han sido observados y almacenados en ordenadores para poder acceder a ellos. Y la observación de hechos, por mucho que los analicemos, no nos conduce a una comprensión general ni a la inteligencia.

En esta verdad ya asumida sobre la investigación científica (y sobre la filosófica) ha surgido un giro moderno: la disponibilidad, más o menos reciente, de cantidades masivas de datos, que al menos en un principio debían fortalecer los sistemas de IA con unos «cerebros» y unas percepciones que antes no se encontraban disponibles. En cierto modo, eso es cierto, pero no en el sentido necesario para escapar a los problemas de la inducción. A continuación, nos fijaremos en el *big data*.

EL FIN DEL *BIG DATA*

El *big data* —o macrodatos— es una idea manifiestamente amorfa que en general se refiere al poder que tienen los conjuntos de datos de gran tamaño para posibilitar análisis y percepciones esenciales por parte de empresas y gobiernos (y de los investigadores en IA). El término apareció impreso por primera vez en un contexto científico en 1997: fue en un documento de la NASA que describía los desafíos que existían para visualizar datos usando la tecnología de gráficos informáticos del momento. Sin embargo, no prendió hasta que se volvió popular durante la década siguiente como un término multifunción en los ámbitos empresarial e informático. Al parecer, el concepto moderno de *big data* afloró en los debates sobre inteligencia empresarial, sobre todo en un informe del grupo Gartner fechado en 2001 y dedicado a los desafíos que presentaba ese campo. El informe destacaba

«tres uves» —volumen, velocidad y variedad— para describir los rasgos de aquellos vastos conjuntos de datos que irían cobrando una importancia cada vez mayor a medida que los recursos informáticos se volvieron más potentes y baratos. Sin embargo, el informe no llegó a utilizar el término de *big data*.⁵ En cualquier caso, el término sí comenzó a aparecer por todas partes a finales de la década de 2000, y en 2014 *Forbes* capturó el bombo y la confusión que lo acompañaban con un artículo titulado «12 definiciones de big data: ¿cuál es la tuya?»⁶

Quizá cueste un poco definirlos con precisión, pero el *big data* —órdenes de magnitud superiores de colecciones de datos— se encuentra en la vanguardia de la revolución informática en la ciencia y en la industria. En 2012, la administración Obama anunció una Iniciativa para la Investigación y el Desarrollo en Big Data con la que se pretendía «solucionar algunos de los desafíos más urgentes a los que se enfrenta el país».⁷ Y al menos una empresa, la firma de análisis empresariales SAS, se apresuró a inventar un nuevo título ejecutivo: vicepresidente de *big data*. Bombo, sin duda, pero la excitación en torno a los macrodatos representaba también un reconocimiento de que, a menudo, una mayor cantidad de datos implicaba una mayor ventaja de cara a analizar los problemas en ordenadores cada vez más potentes.

No obstante, desde un principio se dio una confusión conceptual sobre la manera exacta en que los macrodatos «fortalecían» las percepciones y la inteligencia. Al principio se pensó que los macrodatos mismos eran responsables de la mejora de los resultados, pero, a medida que los métodos de aprendizaje automático fueron tomando vuelo, los investigadores comenzaron a otorgarles el mérito a los algoritmos. El aprendizaje profundo y demás aprendizajes automáticos y las técnicas estadísticas condujeron a mejoras evidentes. Sin embargo, el rendimiento de los algoritmos quedó ligado a los grandes conjuntos de datos. En cualquier caso, la IA estaba demostrando mejoras en su rendimiento, y algunos problemas que hasta entonces no habían tenido solución la encontraron de repente con la entrada de más datos. Y fue esa expansión de la percepción —en los negocios y en la ciencia— lo que los investigadores y los expertos quisieron capturar. En palabras de Jonathan Stuart Ward y Adam Barker, científicos informáticos de la universidad de St. Andrews, «el *big data* está relacionado de manera intrínseca con el análisis de datos y con el descubrimiento de significados derivado de esos datos».⁸ La IA llevaba décadas esforzándose por encontrar

un significado en los datos; de repente, al añadirles más datos aún, ese significado parecía revelarse por todas partes.

En 2013, Viktor Mayer-Schönberger y Kenneth Cukier admitían ya en su *best seller Big Data. La revolución de los datos masivos* que «no existe una definición rigurosa de *big data*», pero sugerían de todos modos que los macrodatos son «la capacidad de la sociedad para emplear la información de una manera novedosa y generar percepciones útiles o bienes y servicios de gran valor», y que su llegada implicaba que en aquel momento había «cosas que se pueden hacer a gran escala, y no desde una escala menor, para extraer nuevas percepciones o crear nuevas formas de valor».⁹ Señalaron historias exitosas en los sectores público y privado que no hubieran sido posibles sin aquel aumento de tamaño en los conjuntos de datos. Piensa, por ejemplo, en la empresa emergente Farecast, fundada en 2004 por Oren Etzioni, emprendedor y profesor de ciencias de la informática de la universidad de Washington, que Microsoft compró en 2008 por más de 110 millones de dólares. Etzioni, que en la actualidad dirige el Instituto Allen de Inteligencia Artificial en Seattle, usó macrodatos en forma de casi doscientos mil millones de registros sobre el precio de los vuelos para encontrar tendencias en sus picos y valles en función de los días previos a la salida. El rendimiento de Farecast subrayó la sensación de que el *big data* significaba nuevos acercamientos y competencias que emergían de aquella gran cantidad de números; partiendo del sistema de referencia de Etzioni, que solo usaba doce mil precios óptimos, el sistema no dejó de mejorar sus predicciones. Al alcanzar los miles de millones de precios óptimos de coste sobre el precio de los vuelos, ofrecía ya un gran valor para el cliente en forma de predicciones certeras sobre el momento en que había que comprar cada billete de avión.

El *big data*, en su momento un término de moda, representa ahora la nueva normalidad para los negocios impulsados por IA de todo el mundo. Walmart creó Walmart Labs para aplicar las técnicas de macrodatos y de extracción de datos a sus desafíos logísticos —comprar, almacenar y enviar la mercadería de manera eficiente en respuesta a la demanda de los consumidores—. Amazon usaba los macrodatos antes de que se pusieran de moda, catalogando y haciendo el seguimiento de las compras en la red, y ahora los usa como datos con los que alimenta los algoritmos de aprendizaje automático que ofrecen recomendaciones de productos, búsquedas mejoradas y otras formas de personalización. Los macrodatos

son una consecuencia inevitable de la ley de Moore: a medida que los ordenadores se vuelven más potentes, las técnicas estadísticas como el aprendizaje automático mejoran y aparecen nuevos modelos de negocio —y todo se debe a los datos y su análisis—. Lo que ahora denominamos ciencia de datos (o, cada vez más, IA) es en realidad una disciplina antigua a la que la ley de Moore y los volúmenes masivos de datos, facilitados en su mayoría por el crecimiento de la red, han dado nuevas alas.

Gobiernos y organizaciones sin fines de lucro no tardaron en sumarse y usar los macrodatos para predecirlo todo, desde la fluidez del tráfico hasta el porcentaje de reincidencia entre los presos con derecho a la libertad condicional. Mayer-Schönberger y Cukier cuentan que Nueva York contrató a unos expertos en *big data* de la universidad de Columbia para que realizaran un modelo predictivo sobre la probabilidad de explosiones de tapas de alcantarilla en la ciudad. (Solo en Manhattan hay más de cincuenta mil tapas de alcantarilla.) El proyecto fue un éxito y se ofreció como ejemplo de la manera en que aquellas nuevas percepciones y competencias resultaban posibles gracias al aumento de escala de los datos. Al fin y al cabo, los obreros humanos no pueden comprobar decenas de miles de tapas de alcantarilla cada día. También otros ámbitos, desde el procesamiento de fichas médicas hasta las iniciativas actuariales del gobierno con respecto al voto y las fuerzas del orden, ofrecen ejemplos que en apariencia apoyan la afirmación según la cual el tamaño y la calidad de los datos —el *big data*— han posibilitado esas nuevas percepciones y competencias.

El éxito de los macrodatos en la industria y otros sectores condujo con rapidez a declaraciones exageradas sobre el poder inferencial de los datos por sí solos. En 2008, las provocativas declaraciones de Chris Anderson, el director de *Wired*, para quien los macrodatos marcarían el fin de la teoría científica, representaron el punto culminante del sensacionalismo relacionado con el tema.¹⁰ Los científicos y demás miembros de la *intelligentsia* se apresuraron a señalar que la teoría es necesaria, cuando menos porque un conjunto de datos no puede pensarse e interpretarse a sí mismo, pero el artículo se mantuvo como una especie de expresión cultural sobre el éxito mareante del diluvio de los datos. En realidad, lo que ocurrió es que, al principio, un batiburrillo de viejas técnicas estadísticas en el uso de la ciencia de datos y el aprendizaje automático en la IA contribuyó a que las expectativas del *big data* emergente se vincularan erróneamente al volumen de datos mismo. Fue una proposición ridícula desde el primer

momento; los datos son hechos y, una vez más, no pueden esclarecer nada por sí solos. Aunque esto solo se haya visto en retrospectiva, los éxitos tempranos del aprendizaje profundo con el reconocimiento de objetos visuales, en los concursos de ImageNet, señalaron el principio de una transferencia de entusiasmo entre el *big data* y los métodos de aprendizaje automático que se benefician de él; en otras palabras, hacia el nuevo y explosivo terreno de la IA.

De modo que el *big data* ha llegado a su cenit, y ahora parece estar desapareciendo del debate popular a la misma velocidad con la que llegaron a él. Tiene sentido que el foco recaiga sobre el aprendizaje profundo porque, al fin y al cabo, son los algoritmos, y no los datos por sí solos, los responsables de haber aplastado a los campeones humanos de go, los que dominan los juegos de Atari, los que conducen coches y todo lo demás. Y, de todos modos, el *big data* ha encontrado un nuevo hogar en la IA contemporánea, mientras que los enfoques basados en datos, como el aprendizaje automático, se benefician de los inmensos volúmenes disponibles de estos para entrenar los modelos y ponerlos a prueba. Tal y como comentaba un observador hace poco, los macrodatos se han convertido en IA *big data*.¹¹

Ya está bien de *big data*. Pero aún tenemos pendiente el tema de la inferencia. Y, en particular, la manera en que los métodos basados en datos, como el aprendizaje automático, pueden superar la brecha entre las simulaciones superficiales basadas en datos y un conocimiento real, obtenido por una capacidad de inferencia más poderosa que la inducción. El problema inmediato es que el aprendizaje automático está inherentemente basado en datos. Ya lo había comentado más arriba, pero a continuación voy a volver sobre ello de manera más precisa.

LA RESTRICCIÓN EMPÍRICA

Los métodos basados en datos sufren en general lo que podríamos denominar como «una restricción empírica». Para entender esa restricción, deberíamos establecer un nuevo aspecto técnico del aprendizaje automático, conocido como «extracción de características». Al abordar un problema concreto, los científicos de IA suelen comenzar identificando los rasgos

sintácticos, o marcas, en unos conjuntos de datos que ayudarán a los algoritmos de aprendizaje a concentrarse en obtener la salida deseada. La ingeniería de rasgos es en esencia una habilidad, y se paga muy bien a los ingenieros y especialistas que disponen de ella de cara a identificar rasgos de utilidad (también a los que tengan el talento para ajustar los parámetros del algoritmo, otro paso del entrenamiento exitoso). Una vez identificados, los rasgos se extraen de manera exclusivamente informática durante las fases de entrenamiento, de prueba y de producción. Y en ese carácter exclusivamente informático radica el quid de la cuestión. Los sistemas de aprendizaje profundo ofrecerían un rendimiento mucho mejor en la difícil tarea del reconocimiento de imágenes si tan solo pudiéramos dibujar una flecha sobre el objeto que deseamos identificar dentro de una foto repleta de objetos y de fondos —usando, pongamos, el *software* de Photoshop—. Pero, por desgracia, ese rasgo de origen humano no se puede añadir a otras fotos ni se puede preparar de esta manera, así que el rasgo no es extraíble sintácticamente y, por tanto, se vuelve inútil. Esa es la semilla del problema. E implica que los rasgos útiles para el aprendizaje automático siempre han de encontrarse entre los datos, y que los seres humanos no podemos proporcionar ninguna pista que la máquina no pueda explotar a su vez «en estado salvaje», al poner a prueba el sistema o después de lanzarlo para su uso.

La extracción de características se realiza en la primera fase, la de entrenamiento, y se vuelve a realizar cuando el modelo ya ha sido entrenado, en lo que se conoce como «la fase de producción». Durante la fase de entrenamiento, se proporcionan datos etiquetados en forma de entrada al algoritmo de aprendizaje. Por ejemplo, si el objetivo consiste en reconocer fotos de caballos, la entrada es la foto de un caballo, y la salida, la etiqueta CABALLO. El sistema de aprendizaje automático («aprendiz») recibe así fotos etiquetadas o marcadas de caballos en forma de parejas de entrada-salida, y la tarea de aprendizaje consiste en simular el etiquetado de imágenes de modo que solo aquellas que muestren un caballo reciban la etiqueta de CABALLO. Se prosigue con el entrenamiento hasta que el aprendizaje genera un modelo —que es un fragmento de código estadístico que representa la probabilidad de un caballo, dada la entrada— que responda (o no) a un criterio de precisión.

Llegado este punto, el modelo generado por el aprendiz se utiliza para etiquetar de manera automática nuevas imágenes, hasta ahora inéditas. Esta

es la fase de producción. A menudo, un bucle de retroalimentación forma parte de esa producción, ya que un humano puede corregir las imágenes de caballos mal etiquetadas y devolvérselas al aprendiz para que continúe capacitándose. Esto puede seguir así de manera indefinida, aunque las mejoras en la precisión se irán volviendo cada vez de menor importancia. La interacción entre los usuarios de Facebook es un ejemplo de bucle de retroalimentación: cuando haces clic sobre un contenido, o etiquetas a un amigo en una foto, estás devolviendo datos al sistema de entrenamiento basado en el aprendizaje profundo de Facebook, que en todo momento estudia y modifica tu analítica de clics para seguir variando o personalizando futuros contenidos.

La restricción empírica es un problema para el aprendizaje automático porque no se puede usar toda la información adicional que desearías suministrarle al aprendiz. A diferencia de las labores de reconocimiento de imagen, que dependen de los datos de píxel como características, en la comprensión del lenguaje hay numerosos problemas que incorporan especificaciones adicionales, esto es, características identificadas por personas que han de ser extraídas por los sistemas cuando se entrenan y utilizan modelos.

Piensa en este problema sencillo del procesamiento de lenguaje, en el reconocimiento de entidades nombradas, donde un conjunto de etiquetas semánticas como PERSONA, ORGANIZACIÓN, PRODUCTO, LOCALIZACIÓN o FECHA son etiquetas de destino o salida, y la entrada es un texto de forma libre, quizá procedente de aportaciones de Facebook. Una empresa o individuo podría querer conocer todas las aportaciones que mencionan a cierta empresa —pongamos, Blue Box, Inc.—. La búsqueda por palabras clave de «Blue Box, Inc.», en la que coincidirían solo esas palabras, podría ignorar referencias más informales, como «Blue Box» o «blue box» o incluso «the box»,^{*b} según el contexto. El objetivo del reconocimiento de entidades nombradas consiste en usar el aprendizaje automático sobre grandes cantidades de aportaciones etiquetadas de modo que esas menciones informales también sean identificadas correctamente como referencias a la empresa. Así, se da la necesidad de una extracción de características: un humano que etiquete todas las menciones a Blue Box, Inc. en una colección de aportaciones que se van a usar en un entrenamiento de datos y que se las mande al sistema, que generará un modelo para etiquetar las menciones a «Blue Box, Inc.» durante la fase de producción.

El sistema de Blue Box depende de las palabras de las aportaciones pero también de sus características: de la presencia en las aportaciones de menciones a empresas. Una vez más, las características son necesariamente sintácticas, porque deben ser extraídas de manera completamente automática durante la fase de producción. Esta es la restricción empírica clave. Las características pueden ser ortográficas, como que se compruebe la presencia de letras mayúsculas, o léxicas, como que se compruebe la presencia de las palabras «*blue*» y «*box*» en ese orden, y pueden incluir información como «que acaben en Inc. o Incorporated». Se puede ejecutar un etiquetador gramatical sobre los datos de entrenamiento para etiquetar categorías como nombres comunes y nombres propios —más características que se detectan sintácticamente—. Sin duda hay otras características posibles. La clave, de nuevo, consiste en que todas ellas, aunque en un primer momento hayan sido identificadas por personas, sean luego extraídas de manera puramente informatizada; de otro modo, el sistema no actuaría automáticamente durante la fase de producción autónoma.

Aquí viene el problema. Algunas señales de menciones a la compañía Blue Box, Inc. requerirán *per se* de una inferencia —por ejemplo, cuando entre los datos aparezcan pronombres y otras referencias—. Esto complica de manera inmediata la labor de aprendizaje. Si leo un estado de Facebook y veo que alguien habla de Blue Box, pero luego se refiere a ella en los comentarios diciendo algo así como «los beneficios de la empresa», la descripción «de la empresa» no es una característica admisible para el sistema de reconocimiento de entidades nombradas. Se ha de encargarse de ello un subsistema de resolución de correferencia, lo cual introduce una tasa de error —hay que tener en cuenta que la correferencia es un problema mucho más difícil que el del reconocimiento de entidades nombradas—. Peor aún, quizá sepamos que Bob está hablando sobre Blue Box, Inc. en un debate sobre el precio de sus acciones, pero, puesto que no encontramos nada al respecto en los datos sometidos a análisis, no existe ninguna característica que pueda ser detectada por el sistema. Alguien podría contar una anécdota sobre el hecho de que el fundador de otra empresa, XYZ, Inc., «adoraba el color azul, y buscaba algo simple y memorable, así que decidió bautizar su sistema operativo como “Blue Box”». Por contexto, aquí «Blue Box» se refiere al producto, no a la compañía, pero el sistema de reconocimiento de entidades nombradas no puede utilizar esa información

contextual durante su entrenamiento. ¿Por qué? Pues porque no puede extraerlo de la mera sintaxis, de su entrada durante la producción.

La restricción empírica forma parte integral del aprendizaje automático. Eso quiere decir que durante la fase de entrenamiento solo se pueden usar características puramente sintácticas que los métodos automáticos descubran entre los datos. Un sistema en verdad inteligente necesita características o señales en un sentido más amplio, que no procedan tan solo de los datos procesados.

Aunque el reconocimiento de entidades nombradas sea una tarea más o menos simple dentro del procesamiento del lenguaje natural, incluso aquí podemos ver las limitaciones inherentes a los enfoques basados exclusivamente en datos. Una mención de Blue Box en una aportación acerca del producto se transforma con facilidad en un falso positivo, y queda etiquetada como si se refiriera a la compañía. Estos ejemplos podrían aparecer en la cola larga de apariciones improbables, pero son bastante comunes en el lenguaje cotidiano, y el aprendizaje automático, limitado por la restricción empírica, no puede abordarlos. Con todo esto pretendo decir que los datos por sí solos, macro o no, y los métodos inductivos como el aprendizaje automático presentan limitaciones inherentes que constituyen barricadas para el avance de la IA. Resulta que el de la inducción representa un problema importante de verdad para la IA moderna. Su ventana hacia el significado está ligada de manera directa a los datos, y eso limita y restringe el aprendizaje.

LA HIPÓTESIS DE FRECUENCIA

Además de la restricción empírica, los métodos de aprendizaje automático dependen de una desafortunada hipótesis de frecuencia. Como sucedía con la restricción empírica, vuelve a tratarse de una consecuencia directa del fundamento enumerativo de la inferencia inductiva —en realidad, se trata de una reafirmación—. Irónicamente, el valor de los macrodatos en el aprendizaje automático es en realidad una muestra de esa suposición: cuanto más, mejor. Los sistemas de aprendizaje automático son solo máquinas contadoras un poco más sofisticadas. Para seguir con el ejemplo de Blue Box, podríamos codificar una lista de características comprobando,

pongamos, si esa palabra o secuencia de dos palabras de algunas aportaciones de Facebook aparece en un diccionario de palabras que incluya nombres de empresas como IBM, Microsoft, Blue Box, etc., o comprobando si le siguen Inc. o LLC, o si se trata de un acrónimo, o si la primera letra está en mayúscula, o si es un nombre común.¹²

La suposición de frecuencia entra en juego porque, en general, a mayor frecuencia de resultados con esa característica, más útil resultará para el entrenamiento. En la ciencia de datos, esto es necesario; si las características de los datos son aleatorias, no se puede aprender nada (como comentamos antes). Pero si existe un patrón, primero, a causa de la restricción empírica, tiene que encontrarse en los datos; y segundo, como consecuencia, la única manera de determinar la intensidad de la asociación entre la entrada y la salida es a través de la frecuencia. ¿Cómo podría haber sido de otro modo? Si cada vez que «Inc.» sigue a un par de palabras la etiqueta en los datos de entrenamiento es EMPRESA, el aprendiz adjudica una probabilidad elevada a «Inc.» como característica de la salida deseada: empresa. Unos patrones que podrían pasar desapercibidos en miles de ejemplos cristalizan cuando los ejemplos son millones. Esa es la hipótesis de la frecuencia.

A las suposiciones de frecuencia se las puede poner patas arriba con la llamada detección de anomalías, que detecta transacciones bancarias fraudulentas o inicios de sesión indebidos. Esos sistemas también dependen de la hipótesis de frecuencia, ya que explotan lo que podríamos denominar la suposición de normalidad. Los sucesos normales hacen que los sucesos anormales resulten más llamativos. Si se pueden agrupar o concentrar miles o millones de inicios de sesión correctos por parte de los empleados, los raritos que quedan fuera del grupo llaman la atención. Por tanto, podrían ser ilegales o indebidos. Una vez más, el aprendizaje automático descubre lo que es normal —y, por tanto, lo que resulta anormal— analizando las frecuencias.

La suposición de frecuencia también sirve para explicar los «filtros burbuja» del contenido personalizado *online*. La persona que desprecia la política de orientación derechista acaba recibiendo solo opiniones de orientación izquierdista y otros contenidos informativos. El sistema basado en el aprendizaje profundo que controla esta salida no hace más que entrenar un modelo que, con el tiempo, pasa a reconocer el patrón de noticias que te gusta. Analiza tus clics y comienza a ofrecerte más de lo

mismo. Las mismas observaciones pueden aplicarse a las sugerencias de Netflix, Spotify, Amazon y otros sitios web que ofrecen búsquedas personalizadas y una experiencia con recomendaciones. Esa conexión entre frecuencia de ejemplo (o características en ejemplos) y aprendizaje automático es intrínseca, esencialmente en el mismo sentido en que inferir que *Todos los cisnes son blancos* se vuelve cada vez más sencillo, porque vas ganando confianza a medida que observas más y más cisnes blancos.

La suposición de frecuencia explica también la dificultad del problema de cola larga de los ejemplos anormales o inesperados. El sarcasmo, por ejemplo, resulta especialmente opaco para el aprendizaje automático, en parte porque es menos frecuente que el sentido literal. Resulta que el recuento funciona bien con ciertas tareas evidentes en la red, pero actúa en contra de las de corte más sutil. Si hay millones de ejemplos de ciudadanos enojados que tuitean «¡Trump es un idiota!», cualquier persona que tuitee «Trump es un idiota» con ánimo sarcástico y como respuesta ingeniosa a los críticos en su intento por superar en astucia a un rival, acabará reducido a una nueva instancia del patrón «Trump-idiota». Para comenzar, el algoritmo de aprendizaje no se encuentra en el negocio del conocimiento, así que para él el ejemplo es solo otra secuencia de palabras. El sarcasmo no es una característica basada en palabras, y tampoco se presenta con la frecuencia del sentido literal. El aprendizaje automático se muestra notablemente obtuso ante esos fenómenos del lenguaje —para disgusto de empresas como Google—. Le encantaría detectar el sarcasmo a la hora de dirigir los anuncios. Por ejemplo, si «¡Tráete la crema solar!» es un comentario sarcástico en una entrada sobre una tormenta de nieve, un sistema de colocación de anuncios sensible al contexto intentaría emitir en su lugar el anuncio de unos calcetines calefactables a pilas.

La hipótesis de frecuencia se vuelve aún más pronunciada cuando la entrada son artículos enteros de prensa que, pongamos, incluyan un texto con la tarea de clasificación mencionada anteriormente. Las noticias «curiosas» o «raras» que abundan en la red como lecturas ligeras representan una pesadilla para el aprendizaje automático, porque el significado de la palabra no es literal. Por ejemplo, un sistema de aprendizaje automático podría clasificar historias que describan sucesos «tontos» pero que técnicamente incluyan referencias a crímenes o delitos como casos genuinos de historias criminales. ¿Por qué no habría de hacerlo? Las noticias peculiares son menos frecuentes que aquellas otras

que proporcionan información de manera directa, así que también aparecen menos en los conjuntos de entrenamiento —y, de todos modos, detectar el motivo por el que resulta peculiar aboca a otro problema con la restricción empírica—. Una noticia se reconoce como tonta o sarcástica cuando alguien comprende la voluntad del autor, lo que pretende comunicar. Pero las palabras en sí de esa historia pueden, por sus frecuencias en los datos de entrenamiento, apuntar a categorías bien definidas como las de política, deportes, crimen, etc. La pieza no podrá ser clasificada, ni comprendida correctamente, a menos que los fragmentos sintácticos que la constituyen —las palabras— se interpreten desde una ventana de significado mucho más amplia. A falta de esta capacidad no inductiva, el sistema de aprendizaje automático pone en valor por defecto las frecuencias, y yerra el tiro. He aquí, por ejemplo, una noticia de Associated Press que apareció publicada una vez en *Yahoo! News*:

¡LOS TACOS O LA VIDA!

Fontana, California — El hambre de carnitas estuvo a punto de provocar una carnicería cuando a un hombre de Fontana le robaron un paquete de tacos a punta de pistola.

El sargento de policía Jeff Decker dijo que la víctima, de 35 años de edad, acababa de comprar tacos por valor de veinte dólares en un puesto callejero el domingo por la noche y regresaba en bicicleta a su casa cuando el sospechoso se plantó ante él y le dijo: «Dame los tacos».

Según Decker, el sospechoso cogió la bolsa con la comida, pegó un puñetazo en la cara a la víctima y se dispuso a huir.

Cuando la víctima le exigió que le devolviera los tacos, el sospechoso le apuntó con lo que pareció ser una pistola y amenazó con matarle antes de escapar.

Un sistema de clasificación de textos identificaría con facilidad esta noticia como una historia criminal: «sospechoso, víctima, huir, pistola». No obstante, la mayoría de los lectores humanos tendrán la impresión de que se trata de una historia cómica —o por lo menos no la vemos como el ejemplo típico de una historia de crímenes—. Se informa sobre los actos criminales porque son graves y preocupantes, pero una frase inicial del tipo «El hambre de carnitas...» señala que la voluntad de Associated Press consiste en informar sobre la historia de manera humorística. Hasta un niño de primaria se dará cuenta de esa voluntad, pero los sistemas de IA estarán encantados de clasificar el artículo como una nueva historia criminal ambientada en Fontana, California. La frecuencia es la asesina del humor. Cuenta el número de artículos periodísticos que incluyen víctimas, pistolas, amenazas y sospechosos que se dan a la fuga. Son de crímenes. El problema del supuesto de frecuencia ante esos ejemplos es que no se dispone de

ninguna solución conocida al usar el aprendizaje automático. El sentido de la historia se pierde dado su método, que analiza las palabras sintácticamente y cuenta la frecuencia de las palabras como señal de una categoría. Incluso con ejemplos más o menos sencillos, como el que nos ocupa, esa senda nos conduce a un callejón sin salida respecto a la inteligencia artificial general.

He aquí otra historia de Associated Press que recogieron diversos periódicos:

Un niño de once años muerde a un *pit bull* para defenderse de su ataque
São Paulo, Brasil — Un niño de once años se ha convertido en el centro de atención de la prensa brasileña después de haber clavado los dientes en el cuello del perro que le había atacado.

Los periódicos locales informaron el jueves de que Gabriel Almeida estaba jugando en el patio de la casa de su tío, en la ciudad de Belo Horizonte, cuando una *pit bull* llamada Tita se lanzó sobre él y le mordió en el brazo izquierdo. Almeida sujetó al perro por el cuello y le devolvió el mordisco... con tanta fuerza que se le rompió un colmillo.

En declaraciones al periódico *O Globo*, Almeida afirmó: «Es mejor perder un diente que la vida».

Unos albañiles que trabajaban cerca del lugar ahuyentaron al perro antes de que pudiera atacar de nuevo.

Sin duda, la noticia cuenta con un lado serio, pero desde luego que no se trata de una historia sobre el ataque de un *pit bull*. Tampoco se trata de la historia sobre un concurso de mordiscos entre un niño brasileño y un perro. Puesto que el niño no sufrió heridas de gravedad, pese a que perdió un diente al morder al perro, queda claro que el motivo para publicar la pieza no fue el de informar sobre el ataque de un can brasileño, sino el de subrayar lo extraño o divertido de ese contraataque sorpresa. Su contenido improbable —niño muerde a perro— es lo que la vuelve noticiable. En un caso como este, la IA y el aprendizaje absoluto no nos ayudan en nada. Nos hacen daño. Yerran el tiro por completo. Los sistemas de aparente inteligencia artificial general que usen solo el aprendizaje automático serán, en el mejor de los casos, unos eruditos idiotas y molestos.

En esencia, la teoría subyacente de la inferencia se encuentra en el meollo del problema. La inducción requiere que la inteligencia sea resultado del análisis de datos, pero la inteligencia llega al análisis de los datos como paso previo y necesario. Siempre cabe esperar que los avances en la extracción de características o el diseño de algoritmos conduzcan en el futuro a una teoría más completa de la inferencia informática, pero

deberíamos mostrarnos muy escépticos a ese respecto. Son precisamente la restricción empírica y la hipótesis de frecuencia las que limitan el alcance y la efectividad de las características detectables —que, al fin y al cabo, se encuentran en los datos para ser analizadas sintácticamente—. Es otra manera de decir lo que los filósofos y científicos de todos los colores aprendieron hace mucho tiempo: no basta con la inducción.

LA SATURACIÓN DE LOS MODELOS

El aprendizaje automático y el *big data* presentan otro problema, conocido como «saturación», que afecta la esperanza de obtener una inteligencia artificial general. La saturación se da cuando, al añadirle más datos —más ejemplos— a un algoritmo de aprendizaje (o técnica estadística), no se suma nada al rendimiento de los sistemas. No existe ningún entrenamiento que pueda prolongarse para siempre e ir ofreciendo una precisión cada vez mayor para un problema. Con el tiempo, el añadido de datos deja de incrementar el rendimiento. Los sistemas exitosos alcanzan una precisión aceptable antes de saturarse; si no es así, el problema no se podrá resolver utilizando el aprendizaje automático. El saturado es un modelo definitivo, que no mejorará por mucho que se le añadan más datos. En algunos casos podría incluso empeorar, aunque las razones para ello son demasiado técnicas para explicarlas aquí.

Rara vez se habla de la saturación del modelo, sobre todo porque numerosos problemas recientes no han dejado de beneficiarse del aumento de datos preparados. Pero los investigadores saben que la saturación resulta inevitable, y que acaba limitando el rendimiento de los sistemas de aprendizaje automático. Allá por 2013, Peter Norvig, director de investigación en Google, confesó a *The Atlantic* la ansiedad que le generaba la saturación: «Podríamos dibujar la siguiente curva: a medida que vamos ganando en datos, ¿hasta qué punto mejora nuestro sistema? —se preguntaba—. Y la respuesta es que continúa mejorando..., pero estamos llegando a un punto en el que obtenemos menos beneficios que en el pasado».¹³

En el momento de escribir estas líneas, el comentario admonitorio de Norvig tiene siete años. Lo más probable es que los concursos de ImageNet

no puedan usar más datos —los mejores sistemas disponen ya del 98 % de precisión (usando la medida de prueba estándar de obtener una etiqueta de destino entre las cinco primeras predicciones de un sistema)—. Pero los coches sin conductor, que creíamos a la vuelta de la esquina, continúan en fase de investigación intensa, y sin duda parte del problema radica en los datos de entrenamiento de las fuentes de vídeo etiquetadas, que no resultan insuficientes en número, pero sí inadecuadas para afrontar problemas de cola larga con escenarios de conducción atípicos que no obstante se han de tener en cuenta por razones de seguridad. Los modelos se están saturando, tal y como predijo Norvig. Es evidente que existe la necesidad de nuevos enfoques. Esas consideraciones son uno de los motivos por los que el llamado crecimiento en escala desde un éxito inicial hasta otro en toda regla resulta ingenuo y simplista. No se puede hacer que un sistema crezca en escala de manera indefinida. El aprendizaje automático —el aprendizaje profundo— no es una fórmula milagrosa.

En sus textos de 1950, Turing se declaró esperanzado ante la idea de que los sistemas computacionales pudieran aprender aquello que desconocían. El aprendizaje automático no era por entonces un término de IA, aunque ya se conocía la posibilidad de crear redes neuronales simples. Pero lo que Turing tenía en mente era una idea ampliada del aprendizaje, más parecida a la versión humana de este. No se podía programar las máquinas con todo el conocimiento que necesitarían, así que debía darse algún tipo de aprendizaje. Pensó que este podría surgir con la inducción. Algunas proposiciones, reflexionó, «vendrían “proporcionadas por la autoridad”, pero otras debería generarlas la máquina misma, por ejemplo, a través de la inducción científica». Desde su localización, en el ecuador del siglo xx, Turing había abandonado sus preocupaciones sobre la necesidad de percepciones externas a los sistemas formales. O, mejor dicho, tenía la esperanza de que encontrarán su lugar en la nueva tecnología computacional.

Sin embargo, ni siquiera los científicos mismos usan la «inducción científica» en el sentido que quiso darle Turing. Realizan conjeturas, acto seguido las ponen a prueba y acto seguido realizan más conjeturas. Turing nunca mencionó el trabajo de Peirce sobre la inferencia lógica. Al parecer, él no tenía un conocimiento sustancial de la inferencia abductiva en el sentido que le dio Peirce.

Y nosotros seguimos buscando sus máquinas capaces de aprender.