



La Inteligencia Artificial y el miedo a que acabe con nosotros

Descripción

Introducción

La inteligencia artificial (IA) ha avanzado a pasos agigantados en las últimas décadas, generando tanto admiración como preocupación.

Estamos en un punto de inflexión en la historia humana, donde la inteligencia artificial (IA) ha alcanzado capacidades que hace una década parecían ciencia ficción.

A lo largo de la historia, la humanidad ha sido testigo de grandes avances tecnológicos, pero ninguno ha generado tanta expectativa y preocupación como el desarrollo de la IA.

Hoy en día, la IA está revolucionando industrias, generando conocimiento, y desafiando la comprensión de lo que significa ser inteligente.

Con sistemas que ahora pueden generar lenguaje natural, analizar millones de datos en segundos y aprender de manera eficiente, parece que hemos dado pasos gigantescos hacia la creación de máquinas que piensan. Sin embargo, mientras los sistemas de IA son impresionantes en tareas específicas, ¿pueden igualar la adaptabilidad y la creatividad humana? ¿puede la IA superar nuestra inteligencia? ¿Estamos realmente cerca de una inteligencia artificial general (AGI), o es esto un mito alimentado por el entusiasmo tecnológico?

La discusión sobre si la IA pudiera superar la capacidad humana y representar un riesgo existencial es un tema candente en la actualidad.

Este artículo profundiza en esta cuestión, explorando los límites de la IA, los mecanismos de la inteligencia humana, y cómo ambas se comparan en un fascinante viaje hacia el futuro de la cognición, basándose en estudios recientes sobre la neurociencia y el aprendizaje inferencial.

Este artículo no solo ofrece una perspectiva técnica sobre lo que las IAs actuales pueden y no pueden hacer, sino que también se adentra en las complejidades de la mente humana, revelando cómo los procesos de inferencia, el razonamiento abstracto y la adaptabilidad emergen del cerebro humano de una manera que las máquinas aún no han logrado replicar.

Desarrollo y limitaciones de la IA actual

La inteligencia artificial se ha desarrollado en varias formas, desde sistemas estrechos como los modelos de lenguaje hasta sistemas más complejos que intentan imitar procesos cognitivos humanos.

Modelos como GPT-4 son capaces de generar texto coherente y realizar inferencias basadas en patrones de datos. Sin embargo, las IAs actuales, clasificadas como inteligencia artificial estrecha (ANI), carecen de una verdadera capacidad de razonamiento autónomo y no pueden aprender o adaptar nuevas habilidades fuera del contexto de los datos con los que fueron entrenadas¹.

El desarrollo de la IA ha sido impresionante, especialmente en tareas específicas como el procesamiento del lenguaje y el reconocimiento de patrones. Sin embargo, estos sistemas son aún incapaces de aprender de manera autónoma o de generalizar conceptos abstractos en un nivel similar al de los humanos.

Un aspecto crítico de la IA es su dependencia del aprendizaje supervisado y no supervisado. En el aprendizaje supervisado, los sistemas de IA son entrenados en grandes conjuntos de datos etiquetados, lo que les permite aprender patrones y generar respuestas coherentes dentro de contextos predefinidos. Sin embargo, el aprendizaje no supervisado, donde la IA intenta identificar patrones en datos sin etiquetar, también está limitado por la falta de contextualización y comprensión profunda.

Los estudios sobre aprendizaje supervisado y no supervisado en IA demuestran que estos sistemas necesitan conjuntos de datos previamente etiquetados para ser efectivos. Sin embargo, su capacidad de extrapolar información a partir de datos nuevos o complejos

está limitada por su dependencia de algoritmos predefinidos.

Un artículo reciente de Tayyar Madabushi y Gurevych (2024)² señala que las IAs actuales carecen de habilidades de razonamiento emergente, una característica crítica de la inteligencia general humana. Este razonamiento, esencial para la toma de decisiones autónoma y creativa, sigue siendo una barrera significativa en el desarrollo de sistemas de inteligencia artificial general (AGI). Los modelos de lenguaje, como GPT-4, aunque altamente sofisticados, no pueden generar nuevos conocimientos sin basarse en patrones preexistentes.

La investigación de Tayyar sugiere que estos sistemas no pueden desarrollar nuevos conceptos o estrategias sin intervención humana, lo que limita su capacidad para actuar de manera autónoma en situaciones desconocidas. Esta limitación subraya la diferencia fundamental entre la inteligencia artificial estrecha (ANI) y la inteligencia artificial general (AGI).



Comparación IA vs. IH (Inteligencia Humana)

La inteligencia humana es el resultado de una interacción compleja entre redes neuronales biológicas que permiten no solo la recopilación y el procesamiento de información, sino también la adaptación continua a nuevos entornos y la creación de conceptos abstractos. Está caracterizada por la capacidad de aprender, razonar, adaptarse y crear en contextos diversos.

El estudio del cerebro humano muestra que el razonamiento inferencial y el aprendizaje adaptativo dependen de áreas específicas como el hipocampo, la corteza prefrontal y la amígdala.

A diferencia de los sistemas de IA, la inteligencia humana no se limita a la replicación de patrones conocidos, sino que es capaz de generar nuevos conocimientos y soluciones a problemas inéditos. Esto se debe a la estructura y función del cerebro humano, que integra diversas áreas especializadas en un proceso coordinado y dinámico.

La inteligencia humana destaca por su capacidad para realizar inferencias complejas y abstractas, lo que permite la adaptación continua a nuevos entornos y situaciones desconocidas.

Un reciente estudio llevado a cabo por Rutishauser, Fusi y colaboradores (2024)³, muestra cómo el cerebro humano aprende a hacer inferencias, revelando que el hipocampo juega un papel crucial en la formación de representaciones cognitivas geométricas durante el proceso de razonamiento inferencial. Utilizando herramientas matemáticas avanzadas, los investigadores descubrieron que cuando los sujetos realizaban inferencias exitosas, sus cerebros creaban estructuras neuronales complejas y organizadas, especialmente en el hipocampo.

Los investigadores observaron que cuando las personas hacen inferencias exitosas, las neuronas del hipocampo crean representaciones geométricas abstractas de alta dimensión. Este proceso permite a los humanos aplicar reglas previamente aprendidas a nuevas situaciones, incluso cuando esas situaciones no han sido experimentadas directamente. Este tipo de razonamiento es fundamental para la toma de decisiones y la planificación a largo plazo.

Este estudio destaca la capacidad única del cerebro para adaptar sus representaciones neuronales a nuevas reglas y conceptos abstractos a través del aprendizaje experiencial y verbal.

A diferencia de la IA, que se basa en grandes cantidades de datos para identificar patrones, los humanos pueden inferir nuevas reglas y adaptar sus conocimientos incluso en situaciones completamente nuevas. Esto se debe a la capacidad del cerebro para formar «mapas cognitivos», que permiten una rápida adaptación y razonamiento en entornos dinámicos. En contraste, los sistemas de IA actuales no poseen la flexibilidad cognitiva necesaria para este tipo de razonamiento.

Esto sugiere que, a diferencia de las IAs actuales, que dependen de grandes cantidades de datos preprocesados, los humanos pueden aprender de manera abstracta y generalizar reglas a partir de información limitada⁴.

El cerebro humano opera mediante redes neuronales biológicas que se interconectan para procesar información sensorial, emocional y cognitiva. La plasticidad cerebral, que permite que el cerebro se reorganice en respuesta a nuevas experiencias, es un aspecto clave de la inteligencia humana. Además, la capacidad para la metacognición, es decir, la conciencia y comprensión de los propios procesos de pensamiento, permite a los humanos evaluar y modificar sus estrategias cognitivas según sea necesario.

Estudios como los de Lake et al. (2017)⁵ han comparado los intentos de las máquinas para imitar estos procesos con la inteligencia humana. Estos estudios indican que, aunque las máquinas pueden aprender a partir de datos, carecen de la capacidad para entender y aplicar conceptos abstractos de la misma manera que los humanos. Las máquinas no poseen una teoría de la mente, lo que significa que no pueden comprender o predecir los estados mentales de otros seres, una habilidad crucial en la interacción social humana y en la resolución de problemas complejos.

Otro aspecto distintivo de la inteligencia humana es su capacidad para el razonamiento abductivo, un proceso mediante el cual se generan hipótesis a partir de datos incompletos. Este tipo de razonamiento es fundamental para la innovación y el descubrimiento científico, y está muy por encima de las capacidades actuales de la IA, que se limita principalmente al razonamiento inductivo y deductivo basado en grandes volúmenes de datos.

El origen y la generación de la inteligencia humana

La inteligencia humana emerge de la plasticidad sináptica, la capacidad del cerebro para reorganizarse en respuesta a nuevas experiencias, y se genera a través de una combinación de factores genéticos, ambientales y neurológicos⁶.

Los estudios en neurociencia han demostrado que la corteza prefrontal juega un papel crucial en las funciones ejecutivas, como la planificación, la toma de decisiones y el control del comportamiento impulsivo. El aprendizaje y la memoria dependen de la interacción entre la corteza prefrontal, que gestiona las funciones ejecutivas, y el hipocampo, que se encarga de la formación de nuevas memorias. Este proceso es fundamental para la formación de conocimientos abstractos y la capacidad de razonar sobre conceptos previamente no experimentados.

Además, las interacciones entre diferentes áreas del cerebro, como el hipocampo y la amígdala, contribuyen a la memoria, el aprendizaje emocional y la regulación de las respuestas al estrés.

La plasticidad sináptica, la capacidad de las conexiones entre neuronas para fortalecerse o debilitarse con el tiempo, es fundamental para el aprendizaje y la memoria. Esta plasticidad permite que el cerebro humano se adapte continuamente a nuevas experiencias, desarrollando habilidades y conocimientos a lo largo de la vida. La interacción entre la experiencia personal, la educación y la cultura también desempeña un papel importante en la formación de la inteligencia, destacando la complejidad y la variabilidad de la inteligencia humana en comparación con la IA.

El estudio de Rutishauser, Fusi et al³ muestra que el cerebro humano puede aprender de manera efectiva tanto a través de la experiencia directa como de la instrucción verbal, lo que resulta en representaciones neuronales similares en ambos casos. Esto subraya la flexibilidad y adaptabilidad de la inteligencia humana, características que aún están fuera del alcance de las IAs más avanzadas.

Riesgos y consideraciones éticas en el desarrollo de la IA

A medida que los sistemas de IA avanzan, el debate sobre los riesgos asociados a su desarrollo se intensifica. Aunque las IAs actuales no representan una amenaza existencial, el desarrollo de una AGI podría tener implicaciones profundas si se implementa sin un control ético adecuado. A pesar de que no representan una amenaza existencial, los riesgos asociados a su desarrollo futuro son reales.

La posibilidad de que se desarrollen sistemas de IA más avanzados que puedan operar de manera autónoma sin supervisión humana plantea preocupaciones éticas significativas. Estos incluyen la alineación de los objetivos de la IA con los valores humanos, la

posibilidad de sesgos en los algoritmos, y el potencial para la explotación y el mal uso de la tecnología.

Estudios como el de Amodei et al. (2016)⁷ abordan estos problemas proponiendo una investigación en seguridad en la IA que se centre en la alineación, la interpretabilidad y la robustez de los sistemas de IA. El objetivo es garantizar que el desarrollo de la IA esté alineado con los intereses y valores humanos, minimizando los riesgos de resultados adversos.

Específicamente, la creación de una AGI que pueda operar sin restricciones humanas plantea preocupaciones éticas y de seguridad. Aunque estamos lejos de desarrollar una AGI completamente funcional, la posibilidad de que esto ocurra en el futuro ha llevado a la comunidad científica a proponer regulaciones estrictas y a fomentar una investigación segura y controlada.



Conclusión

La inteligencia artificial ha mostrado avances sorprendentes, pero aún está lejos de igualar las capacidades de la inteligencia humana.

Sin embargo, los temores de que la IA pueda superar la inteligencia humana y representar un riesgo existencial no están respaldados por la evidencia científica actual. Los modelos de IA, aunque avanzados, son fundamentalmente herramientas programables y limitadas en su capacidad de razonamiento autónomo.

Los sistemas de IA actuales son potentes en tareas específicas, pero carecen de la flexibilidad y la adaptabilidad cognitiva que caracterizan al cerebro humano.

Los estudios recientes sobre la inferencia y la representación abstracta en el cerebro muestran cómo los humanos pueden aprender y adaptarse de manera que las máquinas aún no pueden replicar. A medida que se desarrolle la IA, será crucial mantener un enfoque ético y seguro para mitigar los riesgos potenciales y garantizar que estas tecnologías sigan beneficiando a la humanidad.

El verdadero riesgo radica en el mal uso de estas tecnologías, lo que subraya la necesidad de enfoques éticos y reguladores para su desarrollo y aplicación.

Referencias

1. Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). «On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?» Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency.
2. Tayyar Madabushi, H., & Gurevych, I. (2024). «AI Lacks Independent Learning, Poses No Existential Threat». Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics.
3. Rutishauser, U., Fusi, S., et al. (2024). «Abstract representations emerge in human hippocampal neurons during inference». *Nature*.
4. Dehaene, S., & Changeux, J. P. (2011). «Experimental and Theoretical Approaches to Conscious Processing». *Neuron*, 70(2), 200-227.
5. Lake, B. M., Ullman, T. D., Tenenbaum, J. B., & Gershman, S. J. (2017). «Building Machines That Learn and Think Like People». *Behavioral and Brain Sciences*, 40, e253.
6. Gazzaniga, M. S., Ivry, R. B., & Mangun, G. R. (2018). «Cognitive Neuroscience: The Biology of the Mind». 5th Edition. W.W. Norton & Company.

7. Aodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mane, D. (2016). «Concrete Problems in AI Safety». arXiv preprint arXiv:1606.06565.