



## ¿De qué manera piensa ChatGPT y otros modelos de IA?

### Description

La técnica de "Cadena de Pensamiento" busca entender de qué manera piensa la IA al resolver problemas paso a paso.

### CONTENIDOS

## Inteligencia Artificial Inexplicable: ¿De Qué Manera Piensa la IA?

La inteligencia artificial (IA) ha revolucionado el mundo de la tecnología, y los modelos de lenguaje de gran escala (LLM) son un claro ejemplo de este avance. Sin embargo, [estos sistemas, a pesar de su eficacia, son a menudo inexplicables](#), incluso para sus creadores. La IA inexplicable se refiere a los sistemas de IA cuyas decisiones y procesos internos no pueden ser comprendidos o justificados por los seres humanos. Estos sistemas, como cajas negras, producen resultados útiles, pero el proceso mediante el cual llegan a estas conclusiones no está claro. Los humanos se enfrentan al desafío de comprender y descifrar cómo el algoritmo ha obtenido un resultado determinado.

## Cajas Negras de IA

Las cajas negras en IA se refieren a modelos o sistemas de inteligencia artificial cuyo funcionamiento interno es opaco o difícil de comprender. Aunque [estos modelos pueden producir resultados útiles](#), la manera en que llegan a esas conclusiones no está clara. Un ejemplo de esto es cuando los desarrolladores de Google descubrieron que su IA había adquirido la capacidad de comprender y traducir un idioma originario de Bangladesh sin que hubiese sido específicamente entrenada para ello. Este tipo de comportamientos resalta la necesidad de [entender mejor los procesos internos de la IA](#).



Algunos investigadores utilizan técnicas de la psicología humana para estudiar los LLM, tratándolos como sujetos de estudio. Esto ha revelado comportamientos sofisticados que emergen de circuitos subyacentes simples, proporcionando una nueva perspectiva sobre la IA.

## Herramientas XAI Para Resolver de qué Manera Piensa la IA

Para abordar la complejidad de los LLM, los investigadores han desarrollado herramientas de IA explicable (XAI) para intentar entender cómo funcionan estos modelos. La IA explicable es un conjunto de procesos y métodos que permiten a los usuarios humanos comprender y confiar en los resultados creados por los algoritmos de aprendizaje automático o machine learning (ML). La IA explicable se utiliza para describir un modelo de IA, su impacto previsto y sus posibles sesgos. Aunque se han logrado avances, la XAI sigue siendo un campo en desarrollo, buscando formas

---

de hacer que la IA sea más transparente y comprensible.

Te Puede Interesar:

## Comportamiento Errático de los LLM

Los Modelos de Lenguaje de Gran Escala (LLM) son algoritmos avanzados de aprendizaje profundo que pueden realizar una amplia gama de tareas relacionadas con el procesamiento del lenguaje natural (NLP). Sin embargo, [a pesar de su eficacia, los LLM pueden comportarse de manera impredecible](#). Este comportamiento errático puede ser evidente en situaciones donde el modelo genera respuestas que no se alinean con las expectativas del usuario o las normas sociales aceptadas. Por ejemplo, un LLM puede generar respuestas que parecen arbitrarias o incluso engañosas. Este comportamiento errático plantea desafíos en la comprensión y aplicación segura de los LLM. Es importante desarrollar métodos para entender y mitigar estos comportamientos erráticos para garantizar que los LLM se utilicen de manera efectiva y segura.



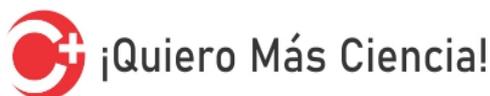
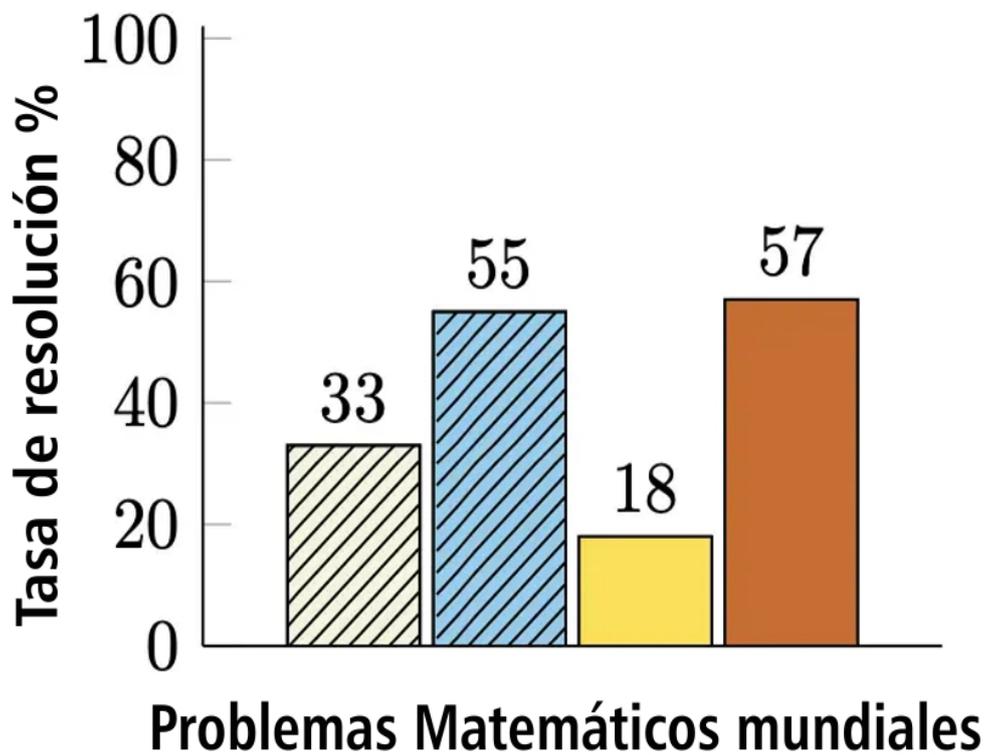
---

a. Cuando las máquinas pueden reconocer y responder adecuadamente a nuestras emociones, la conexión humano-máquina se fortalece. La implementación de interfaces emocionales en la IA, como la capacidad de los chatbots para reconocer el tono emocional en el lenguaje, marca un paso crucial hacia la creación de sistemas más intuitivos y centrados en el usuario.

## Técnica de "Cadena de Pensamiento"

La técnica de "Cadena de Pensamiento" es un método recientemente desarrollado que alienta a los LLM a explicar su razonamiento. Esta técnica se basa en la idea de que al mostrarle al LLM algunos ejemplos donde se explica el proceso de razonamiento, el LLM también mostrará el proceso de razonamiento al responder al prompt. Esta explicación del razonamiento a menudo conduce a resultados más precisos. Sin embargo, es importante destacar que, según Wei et al., "CoT solo produce mejoras de rendimiento cuando se usa con modelos de alrededor de 100 mil millones de parámetros". Los modelos más pequeños escribieron cadenas de pensamiento ilógicas, lo que condujo a una precisión peor que la del prompting estándar.

-  **GPT-3 Afinada**
-  **Mejor prioridad**
-  **Prompt estandard PaLM**
-  **Prompt de cadena de pensamiento PaLM**



Comparación de modelos en la prueba de referencia GSM8K (Wei et al.)

## De Qué Manera Piensa la IA: Escaneos Cerebrales y la Honestidad

La inteligencia artificial (IA) está emergiendo como [una herramienta en la lucha contra enfermedades como el Alzheimer](#). Un grupo de investigadores, liderado por Paul Thompson de la Universidad del Sur de California en Los Ángeles, ha diseñado algoritmos capaces de examinar miles de imágenes cerebrales. Estos avances no solo abren nuevos horizontes en la comprensión de la enfermedad, sino que también prometen transformar radicalmente el enfoque hacia su diagnóstico, superando los obstáculos tradicionales y desbloqueando nuevas posibilidades para combatir este enigmático trastorno. Al medir la actividad neuronal, han creado representaciones

---

matemáticas de la honestidad, abriendo la posibilidad de detectar deshonestidad en tiempo real. Este tipo de comportamientos resalta la necesidad de entender mejor los procesos internos de la IA.

## Edición de Redes Neuronales

Las redes neuronales artificiales (ANN) están formadas por capas de nodos, que contienen una capa de entrada, una o varias capas ocultas y una capa de salida. Cada nodo, o neurona artificial, se conecta a otro y tiene un peso y un umbral asociados. Si la salida de un nodo individual está por encima del valor de umbral especificado, dicho nodo se activa y envía datos a la siguiente capa de la red. De lo contrario, no se pasan datos a la siguiente capa de la red. Las redes neuronales se basan en entrenar datos para aprender y mejorar su precisión con el tiempo. Técnicas avanzadas permiten a los científicos editar las redes neuronales de la IA, cambiando respuestas específicas sin necesidad de reentrenar todo el modelo. Esto podría ser clave para actualizar y corregir conocimientos en la IA.

## Para seguir pensando

La [inteligencia artificial](#) ha avanzado a pasos agigantados, pero su funcionamiento interno sigue siendo un misterio para muchos. Los modelos de lenguaje de gran escala (LLM) son un claro ejemplo de este progreso, pero su comportamiento puede ser impredecible. Para abordar estos desafíos, se han desarrollado herramientas de IA explicable (XAI) que buscan hacer que la IA sea más transparente y comprensible. A medida que la IA avanza, los humanos se enfrentan al desafío de comprender y descifrar cómo el algoritmo ha obtenido un resultado determinado. La psicología de las máquinas y la técnica de "Cadena de Pensamiento" son algunos de los enfoques que se están explorando para entender mejor estos sistemas. Sin embargo, aún queda mucho por descubrir en este emocionante campo de estudio.