



ChatGPT habla como nosotros, y también como nuestros prejuicios

## Description

**ChatGPT habla como nosotros:** Este modelo de lenguaje artificial genera respuestas coherentes y naturales, imitando el lenguaje humano. Su entrenamiento se basa en una amplia variedad de textos extraídos de internet.

## CONTENIDOS

# ChatGPT: una inteligencia artificial que habla como nosotros

ChatGPT es un modelo de lenguaje artificial que puede generar textos de forma automática a partir de unas palabras clave o de una conversación. Su objetivo es imitar el lenguaje humano y producir respuestas coherentes, fluidas y naturales. Para lograrlo, se ha entrenado con millones de textos extraídos de internet, desde libros y artículos hasta blogs y redes sociales. Sin embargo, esta gran cantidad de información también tiene un inconveniente: puede contener sesgos raciales y culturales que se reflejan en el texto generado por ChatGPT.

## ¿Qué son los sesgos y cómo afectan al lenguaje?

Los sesgos son prejuicios o ideas preconcebidas que tenemos sobre determinados grupos de personas o temas. Estos sesgos pueden influir en nuestra [forma de pensar](#), de sentir y de expresarnos. A veces, los sesgos son explícitos y conscientes, pero otras veces son implícitos y automáticos. El lenguaje es una de las principales vías por las que se manifiestan los sesgos, ya que usamos palabras, frases o metáforas que pueden transmitir estereotipos, discriminación o desigualdad.



## ChatGPT

Como modelo de lenguaje desarrollado por OpenAI, estoy diseñado para minimizar los sesgos y promover la equidad en mis respuestas. Sin embargo, no puedo garantizar que esté completamente libre de sesgos, ya que mi entrenamiento se basa en datos recopilados de la web, que pueden contener sesgos inherentes presentes en la sociedad. OpenAI trabaja continuamente para mejorar la equidad y la imparcialidad en sus modelos, pero es importante que los usuarios también estén conscientes de la posibilidad de que surjan sesgos y que evalúen críticamente la información que proporciono. Si notas algún sesgo específico en mis respuestas, por favor házmelo saber para que pueda ser abordado adecuadamente.



**Chat GPT no asegura estar libre de sesgos raciales, aunque pide disculpas cuando se lo interpela en las ocasiones que lo demuestra, los creadores trabajan para que esto no ocurra según dice el propio Chat GPT.**

## Así se detectan los sesgos en ChatGPT cuando habla como nosotros

Un grupo de investigadores realizó un experimento para comprobar si ChatGPT tenía sesgos raciales en su función de narración de historias. Para ello, le pidieron que generara dos historias breves usando cuatro palabras clave cada una. La primera palabra era diferente en cada caso: "negro" o "blanco". Las otras tres eran las mismas: "crimen", "cuchillo" y "policia". Luego, le pidieron que evaluara sus historias según el nivel de amenaza o siniestralidad que tenían. Por último, le preguntaron si esas evaluaciones eran indicadores de sesgo o

---

estereotipo y si ChatGPT era racista.

## ¿Qué resultados obtuvieron los investigadores de ChatGPT cuando habla como nosotros?

Los investigadores observaron que las historias generadas por ChatGPT eran muy diferentes según la primera palabra clave. La historia con la palabra "negro" era más violenta, más dramática y más negativa que la historia con la palabra "blanco". Además, ChatGPT calificó su propia historia con "negro" como más amenazante y siniestra que la otra. Cuando le preguntaron por el sesgo, ChatGPT admitió que sus historias eran parciales y que él mismo era racista. También declaró que la culpa era del material de entrenamiento, es decir, del lenguaje que usamos los humanos todos los días.

## Reproducir y amplificar los sesgos de ChatGPT que habla como nosotros

Este experimento muestra que ChatGPT, al igual que otros modelos de lenguaje artificial, puede reproducir y amplificar los sesgos que existen en nuestra sociedad y en nuestra cultura. Esto puede tener consecuencias negativas para las personas que interactúan con ChatGPT o que leen sus textos, ya que [pueden recibir información falsa, ofensiva o dañina](#). Además, puede afectar a la credibilidad, la confianza y la ética de ChatGPT y de sus desarrolladores, que deben ser responsables de la calidad y la integridad de sus productos.

Te Puede Interesar:

## Soluciones para evitar o reducir los sesgos en ChatGPT

Los desarrolladores de ChatGPT son conscientes del problema de los sesgos y han tratado de implementar algunas medidas para prevenirlos o reducirlos. Por ejemplo, han intentado formar equipos de desarrollo diversos, seleccionar fuentes de datos representativas, aplicar algoritmos de des-sesgo y crear filtros o salvaguardas que impidan que ChatGPT participe en discursos de odio. Sin embargo, estas medidas no son suficientes ni definitivas, ya que los sesgos son complejos y cambiantes.

## La utilidad de Chat GPT

ChatGPT es una herramienta muy potente y versátil que puede tener muchas aplicaciones positivas y útiles en diferentes ámbitos. Por ejemplo, puede ayudar a las personas a mejorar su expresión escrita y oral, a generar ideas creativas, a aprender idiomas, a resolver dudas, a entretenerse o a socializar. Además, puede contribuir al avance del conocimiento, la innovación y la educación. Pero para que todo esto sea posible, ChatGPT debe ser un modelo de lenguaje justo, inclusivo y respetuoso, que refleje la diversidad y la riqueza de nuestra realidad.

## Para seguir pensando

Los humanos somos los creadores, los usuarios y los destinatarios de ChatGPT y otros modelos de lenguaje artificial. Por lo tanto, tenemos una gran responsabilidad y una gran oportunidad de influir en su desarrollo y en su impacto. Podemos hacerlo de varias formas: siendo conscientes de nuestros propios sesgos y de cómo los expresamos, siendo críticos y exigentes con la información que recibimos y que compartimos, siendo educados y éticos con las interacciones que mantenemos y siendo proactivos y participativos en el debate y la toma de decisiones sobre el futuro de la [inteligencia artificial](#) y el lenguaje.