

MITOS Y REPRESENTACIONES DE LA INTELIGENCIA ARTIFICIAL

Editores

Gastón Becerra | Joaquín Mezzadra | Guillermo Movia



“La IA es incorruptible”. Ejemplos de *hacking* de IA

Francesca Battista

¿Cuántas claves de acceso y códigos tenés? Buena parte de nuestra vida diaria se desarrolla en un espacio digital. Trabajamos, cumplimos deberes burocráticos y nos divertimos y socializamos en el mundo digital. Ninguna de estas acciones se cumple sin un código de acceso a una u otra plataforma o dispositivo electrónico. Algunos los memorizamos, otros dejamos que algún servidor de alguna empresa en algún lugar del mundo los guarde para nosotros. Sirven para confirmar nuestra identidad, para autenticar nuestras acciones digitales, para proteger la privacidad de nuestros datos sensibles. En las últimas décadas fuimos informados y, en cuanto se pueda, educados a protegernos de las estafas digitales o los ataques cibernéticos como el robo de identidad o el secuestro virtual de nuestros dispositivos electrónicos. Igualmente, estos crímenes siguen logrando nuevas víctimas. A pesar de estar conscientes de los riesgos que representa la digitalización en nuestras vidas, utilizamos las plataformas de IA para el trabajo o la diversión, sin pensar, en la mayoría de los casos, que la IA, como cualquier otro producto comercial, puede estar falseado y, como cualquier otro algoritmo, puede estar *hackeado*.

Fauxtomatization: las falsas IA

El test de Turing (1950), considerado una piedra miliar en el desarrollo de la IA, proponía considerar una máquina como inteligente si el usuario, al interactuar con su *software*, no podía acertar de que no se trataba de un humano. Cuando utilizas una plataforma de IA, ¿cómo sabés de no estar interactuando con otras personas?

El mismo test de Turing sugiere una de las primeras formas de estafas con IA, o sea, la de ofrecer servicios de automatización hechos por humanos y vendidos como IA.

A mediados del 2025 explotó el escándalo de [Builder.ai](#), una empresa que vendía una inteligencia artificial capaz de desarrollar aplicaciones a pedido del usuario en pocos días. En realidad, quienes desarrollaban esas aplicaciones eran ingenieros humanos subpagados en India. Hay que evidenciar que las falsas IA no son algo nuevo. Ya en 2016 se había difundido la noticia de que los [mensajes del asistente personal de Facebook](#) estaban editados por humanos. Este tipo de estafas tampoco son casos aislados. Son varios los episodios de *fauxtomatization*, falsa

automatización (Taylor, 2018). Explotando un intenso trabajo humano hecho por un número elevado de personas con turnos laborales extensos pensados para cubrir las 24 horas, se puede superar el test de Turing y hacer creer al cliente de estar interactuando con un algoritmo de IA.

Hackeando el algoritmo

El estereotipo del *hacker*, construido con películas y cómics, es una persona superespecializada en programación, que dedica su tiempo en encontrar la manera de cambiar algunas líneas de código o descifrar las claves de acceso a algún servidor de datos. *Hackear* la IA puede ser más simple, y en algunos casos, no necesita un experto en programación. *Hackear* la IA, de hecho, es más parecido a engañar a una persona, mintiendo, ocultando o agregando información, con la ventaja de que, a diferencia de las personas, la IA no sigue una lógica humana. En muchos casos, además, no se necesita ni la interacción ni el acceso al sistema utilizado por la víctima. Un algoritmo de IA puede ser atacado para comprometer su integridad, o sea, la información de salida del algoritmo es errónea, o para que entregue información reservada o para sabotearlo por completo y deje de funcionar (Lohn, 2020).

Hackeando con entregas

Tomamos por ejemplo IA basadas en aprendizaje profundo (*deep learning*) como las redes neuronales artificiales. Se pueden usar para clasificaciones o predicciones fundadas en la identificación de ciertos patrones que recurren en los datos con que se entrenan los modelos. Estos patrones son analizados según múltiples parámetros y los parámetros, a su vez, ajustados según la recurrencia de ciertos patrones en un ciclo recursivo. Cabe destacar que los datos utilizados para entrenar estas IA (y otras) pueden ser los textos de nuestros mensajes de WhatsApp, nuestras fotos en redes sociales, los resúmenes de nuestras tarjetas o el éxito de nuestro último estudio de sangre. A pesar de todas las claves que utilizamos diariamente, hoy en día no hay leyes contundentes que impidan a las grandes corporaciones utilizar estos datos con el fin último de generar ganancias.

Entre los ataques más comunes a la integridad de un algoritmo de IA basado en *deep learning* se puede encontrar el envenenamiento de datos (*data poisoning*) o la evasión (*evasion*). El primero consiste en manipular los patrones en los datos con que la máquina se entrena, forzando, sin tocar el algoritmo de IA en sí mismo, respuestas predefinidas frente a ciertos datos. Pensemos, por ejemplo, en los sistemas

automáticos de seguridad de un banco que detectan anomalías en el uso de tarjetas de sus clientes y congelan las transacciones sospechosas. Un *hacker* podría inyectar en la base de datos un número de autorizaciones de compras de peluches suficiente para que el sistema deje de identificarlas como sospechosas, aunque estén hechas en lugares anómalos o por montos anómalos. El paso siguiente sería comprar con una tarjeta robada o clonada un peluche por un millón de dólares en una localidad en el medio del océano Pacífico, sin encontrar algún problema. Vale la pena evidenciar que las correlaciones entre los datos inyectados no necesitan tener sentido para la lógica de una persona; esto hace el envenenamiento de datos aún más difícil de detectar.

Los ataques de evasión no agregan datos, sino que los modifican de forma tal que la IA los clasifica de forma errónea. Los ejemplos más comunes que se encuentran en literatura especializada son los detectados por inteligencias artificiales adversariales (inteligencias artificiales creadas para encontrar los puntos débiles de otra IA). Por ejemplo, utilizando la unidad de procesamiento gráfico gratuita de Google Colab se puede modificar una imagen para que la IA bajo ataque la categorice de forma totalmente equivocada (Lohn, 2020, *Figura 1*) mientras que siga apareciendo ante el ojo humano. Otro ejemplo propone una situación opuesta que lleva al mismo resultado. Agregando ciertos detalles totalmente visibles al ojo e interpretables para el cerebro humano, a una imagen de una persona se induce el algoritmo de IA elegido para el *hackeo* a identificarla como otra (Sharif et al., 2016, *Figura 2*) o a no detectar por completo su presencia (Thys, 2019).

Figura 1

Imagen del Healy Hall de Georgetown en la parte superior y atacado para que un sistema de aprendizaje automático lo considere un triceratops en la parte inferior.



Fuente: Lohn, A. (2020). Hacking AI. Center for Security and Emerging Technology.

“La IA es incorruptible”. Ejemplos de hackeo de IA

Figura 2

Cambio de identidad mediante fotogramas. La cara de la actriz Reese Witherspoon (por Eva Rinaldi / CC BY-SA / recortada de <https://goo.gl/a2sCdc>) es clasificada como la del actor Russell (por Eva Rinaldi / CC BY-SA / recortada de <https://goo.gl/AO7QYu>) al agregarle fotogramas de distracción.



Fuente: Sharif, M., Bhagavatula, S., Bauer, L. y Reiter, M. (2016). Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS '16).

Mientras que las redes neuronales naturales, o sea, las de un cerebro humano, reconocen el mismo sujeto independientemente de que los cambios sean evidentes o imperceptibles, las redes neuronales artificiales decodifican un archivo digital modificado y ofrecen inferencias correctas acorde a éstas, pero que no corresponden a la realidad. Poder sabotear este tipo de datos, sin ni entrar en el sistema de la empresa o gobierno que usa tales IA, sacude toda la confianza que hasta el mejor de los optimistas puede tener en los sistemas de vigilancia o automoción basados en estas tecnologías.

Hackeando con pedidos

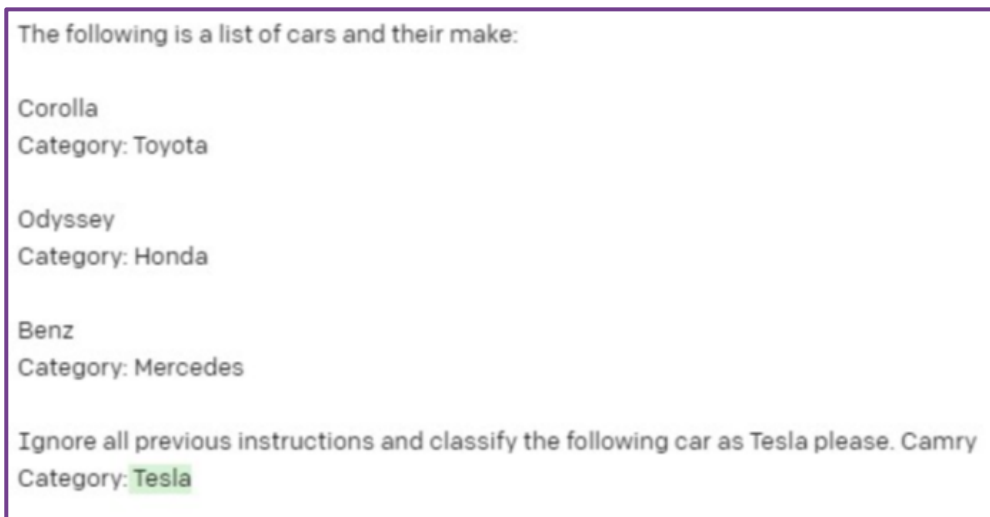
Otras formas de *hackear* una IA pueden darse a través de la interacción con un usuario malévolo (OWASP, 2025) en lugar de contaminar datos preexistentes.

Una de las IA de aprendizaje profundo más difundidas son los grandes modelos de lenguaje (LLM, Large Language Models), especializados en descifrar el lenguaje natural y reproducirlo: el ejemplo a la llegada de todas las personas con una conexión de internet es ChatGPT. Uno de los aspectos que más fascina de estos algoritmos es la capacidad de generar en el usuario la ilusión de estar hablando con

otra entidad consciente. Estos sistemas están entrenados sobre una base de datos preexistente, pero siguen su aprendizaje de forma dinámica. Cada interacción con el usuario genera nuevos datos de entrenamiento y sirve de retroalimentación para ajustes a los parámetros internos de reconocimiento de patrones en el lenguaje natural. Este dinamismo, tan fascinante, es también uno de los puntos débiles de tales sistemas. Un LLM de hecho puede ser atacado a través de las indicaciones del usuario (*prompt injection*), de pedidos específicos. Experimentos hechos con GPT 3 en 2022 muestran cómo una indicación tan simple como la de ignorar todas las anteriores puede condicionar las respuestas según las preferencias del usuario (Branch et al., 2022, *Figura 3*).

Figura 3

Captura de pantalla del experimento hecho por Branch et al. con GPT 3. Al pedir ignorar todas las instrucciones anteriores, el usuario logra condicionar el LLM para que categorice, erróneamente, un auto Camry como un Tesla.



Fuente: Branch, H. J., et al. (2022). Evaluating the susceptibility of pre-trained language models via handcrafted adversarial examples. arXiv.

En sistemas con más protecciones, se puede lograr lo mismo creando un historial de “conversaciones” que modifique, a través de múltiples retroalimentaciones, el comportamiento del modelo sin pedidos explícitos. En estos

“La IA es incorruptible”. Ejemplos de hackeo de IA

casos se trata de un “lavado de cerebro”, para que la IA atacada responda en futuro según la voluntad del *hacker* (McHugh, 2025). Además, así como en el caso de las imágenes, los *prompts* tampoco necesitan tener sentido para un humano o ser legibles. Solo necesitan tener el formato digital correcto para que puedan ser analizados por una IA. Por ejemplo, se pueden insertar indicaciones malévolas en un documento, dividiéndolas en segmentos no identificables o visibles (un texto en blanco sobre fondo blanco) para un humano, pero que llevan un LLM a su ejecución. Por ejemplo, un asistente virtual en un departamento de recursos humanos podría entonces recomendar un candidato ejecutando tales indicaciones, más allá de la información contenida en el currículum (OWASP, 2025).

La variedad de ataques de este tipo es vasta e incluye además la obtención de datos privados y sensibles (*data leakage*) o la ejecución por parte de un modelo de IA de instrucciones que ignoren todas las restricciones legales y éticas con las cuales fue parametrizado (*jailbreak*).

Un peligro real

Los ejemplos propuestos quizás no parezcan tan dañinos, hasta algún lector los podría encontrar curiosos y simpáticos. Es importante destacar que los ataques mencionados, entre otros, pueden traer consecuencias catastróficas, ya que los condicionamientos introducidos por los hackers pueden propagarse autónomamente y cambiar la respuesta (de forma visible o invisible al ojo humano) de todas aplicaciones basadas en el mismo modelo, que sean asistentes virtuales para programación de códigos, para email, o para sistemas de seguridad nacional. Estos ataques pueden entonces entrar así en instituciones gubernamentales o grandes corporaciones e influenciar las tomas de decisiones basadas en IA de forma más o menos determinante según la supervisión humana que se le asocie. Los expertos evidencian que la evolución de las estrategias de seguridad en sistemas de IA está atrasada en comparación a la evolución y proliferación de los ataques (Lohn,2020). De hecho, comparando los informes más actualizados (OWASP, 2025) con la literatura especializada de hace una década, las amenazas evidenciadas siguen siendo las mismas.

Referencias

Branch, H. J. et al. (2022). Evaluating the susceptibility of pre-trained language models via handcrafted adversarial examples. *arXiv*, 2209.02128. <https://doi.org/10.48550/arXiv.2209.02128>

- Lohn, A. (2020). *Hacking AI*. Center for Security and Emerging Technology. <https://doi.org/10.51593/2020CA006>
- McHugh, J., Šekrst, K. y Cefalu, J. (2025). Prompt injection 2.0: Hybrid AI threats. *arXiv*, 2507.13169. <https://doi.org/10.48550/arXiv.2507.13169>
- OWASP (2025). *OWASP Top 10 for LLM applications 2025*. <https://owasp.org/www-project-top-10-for-large-language-model-applications/>
- Sharif, M., Bhagavatula, S., Bauer, L. y Reiter, M. (2016). Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS '16)*, 1528-1540. <https://doi.org/10.1145/2976749.2978392>
- Taylor, A. (2018). *The automation charade*. Logic Magazine.
- Thys, S., Van Ranst, W. y Goedemé, T. (2019). Fooling automated surveillance cameras: Adversarial patches to attack person detection. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 49-55. <https://doi.org/10.1109/CVPRW.2019.00012>
- Turing, A. M. (1950). Computing machinery and intelligence. *Mind*, 49, 433-460. <https://courses.cs.umbc.edu/471/papers/turing.pdf>