

Resumen ejecutivo

“La inteligencia artificial no es un sustituto de la inteligencia humana; es una herramienta para amplificar la creatividad y el ingenio humanos”.

Fei-Fei Li²²



LLM: contexto, definición y regulación

1. La inteligencia artificial generativa (GenAI), y dentro de ella los modelos de lenguaje a gran escala (LLM) representan un avance significativo en el campo de la IA, que define una nueva generación de interfaz hombre-máquina en la que la comunicación se realiza mediante el lenguaje natural, y con aplicaciones revolucionarias en todos los sectores, incluyendo la educación, la salud, las finanzas y el comercio. Sin embargo, su desarrollo y uso también conllevan riesgos y desafíos importantes que deben abordarse.
2. Los LLM son modelos de IA entrenados para reconocer, generar, traducir y resumir grandes cantidades de texto. Utilizan arquitecturas como los *transformers* y se entrenan con vastos conjuntos de datos para aprender patrones y estructuras lingüísticas. Su eficacia depende del tamaño en términos de número de parámetros, la estructura, la diversidad de los datos de entrenamiento y la sofisticación de sus algoritmos.
3. La evolución de los LLM ha sido muy rápida, desde los primeros modelos basados en reglas hasta los actuales basados en *transformers*. Hitos importantes incluyen la introducción de la arquitectura *transformer* y los mecanismos de autoatención, y los primeros LLM comerciales, como GPT. El año 2023 fue clave, con una mayor accesibilidad, contribuciones globales y la proliferación de los LLM de código abierto.
4. Los LLM tienen numerosas aplicaciones, como la creación y mejora de contenido, el análisis y organización de información, y la interacción y automatización de tareas. Con la emergencia de LLM multimodales, se están abriendo nuevas posibilidades en la generación de contenido audiovisual y experiencias interactivas enriquecidas.

5. Los reguladores están tomando medidas para abordar los riesgos y oportunidades de la IA, con iniciativas como el AI Act de la UE, el AI Bill of Rights de EE.UU. y la Declaración de Bletchley. Algunos requisitos clave incluyen transparencia, privacidad, equidad, seguridad, responsabilidad y supervisión humana.

Desarrollo y despliegue de LLM

6. El desarrollo de LLM implica varios componentes y decisiones críticas, como la selección y preprocesamiento de datos, la *tokenización* y los *embeddings*, el preentrenamiento, la *cuantización* y el *fine-tuning*. En particular, el elevado coste del entrenamiento suele derivar en la elección de usar un modelo preentrenado o un modelo de código abierto, y limitarse a hacer *fine-tuning* con datos relativos a la aplicación que se quiere desarrollar. La implementación requiere consideraciones de integración, monitoreo y aspectos éticos y legales.
7. El entrenamiento de los modelos es un aspecto crucial que influye en su eficacia. Factores como la cantidad y calidad de los datos de entrenamiento, la arquitectura del modelo y los algoritmos de aprendizaje utilizados pueden tener un impacto significativo en el rendimiento y la generalización de un LLM.
8. La arquitectura más común para los LLM son los *transformers*, que utilizan mecanismos de autoatención que permiten al modelo encontrar relaciones entre distintas partes del texto, procesarlo y generar nuevo texto. Han demostrado un rendimiento excepcional en diversas tareas de procesamiento de lenguaje natural. Variantes y extensiones buscan mejorar su eficiencia y escalabilidad.

²²Fei-Fei Li (n. 1976). Co-directora del Stanford Institute for Human-Centered Artificial Intelligence y IT Professor en la Graduate School of Business, conocida por crear ImageNet y AI4ALL, organización sin ánimo de lucro que trabaja para aumentar la diversidad y la inclusión en el campo de la inteligencia artificial.

9. LLMOps es una metodología para gestionar el ciclo de vida completo de los LLM, abordando desafíos como la gestión de grandes volúmenes de datos, el escalado de recursos computacionales²³, la monitorización y el mantenimiento, el versionado y la reproducibilidad.
10. Los principales retos de los LLM incluyen sesgos y alucinaciones, falta de explicabilidad y transparencia, calidad y accesibilidad de los datos, problemas de privacidad y seguridad, y alto consumo de recursos. También existen desafíos de dependencia, riesgos de uso malicioso, cuestiones de propiedad intelectual y escalabilidad.

Marco de validación de LLM

11. La validación de los LLM es crucial para garantizar su uso seguro y responsable, y conviene adoptar una perspectiva amplia que abarque los diversos riesgos asociados. Un marco de validación multidimensional debe cubrir aspectos como el riesgo de modelo, la gestión de datos, la ciberseguridad, los riesgos legales y operativos, la ética y la reputación.
12. La validación de LLM debe articularse mediante una combinación de métricas cuantitativas y técnicas de evaluación humana. La selección de técnicas dependerá de las características del caso de uso, como el nivel de riesgo, la exposición pública, el procesamiento de datos personales y la línea de negocio.
13. Las tendencias emergentes en la validación de LLM incluyen la explicabilidad²⁴, el uso de LLM para explicar otros LLM, puntuaciones por atribución, validación continua, enfoques colaborativos, ingeniería de *prompts*, alineación ética y regulatoria, y técnicas de desaprendizaje (*machine unlearning*).

Caso práctico

14. El caso práctico presentado ilustra la aplicación de un marco de validación personalizado a un *chatbot* de políticas internas de una compañía. El proceso abarcó la definición del caso, el diseño del enfoque de validación, la ejecución de pruebas cuantitativas y cualitativas, y la interpretación de resultados.
15. Los resultados de la validación del *chatbot* mostraron un desempeño general satisfactorio, con fortalezas en precisión, consistencia, adaptabilidad y escalabilidad. Se identificaron áreas de mejora en explicabilidad, mitigación de sesgos y seguridad. Se recomendó proceder con la implementación, aplicando las mejoras sugeridas y estableciendo un plan de monitoreo y perfeccionamiento continuo.

Conclusión

16. En conclusión, los LLM tienen un potencial significativo para transformar diversos sectores, pero su desarrollo y despliegue también conllevan retos significativos en áreas como la transparencia, la equidad, la privacidad y la seguridad. Para aprovechar los beneficios de los LLM de manera responsable, es fundamental establecer un marco sólido de gobierno de la IA que aborde estos desafíos de manera integral, incluyendo un enfoque riguroso y multidimensional de validación que cubra todo el ciclo de vida de los modelos. Solo así se podrá garantizar que los LLM sean fiables, éticos y estén alineados con los valores y objetivos de las organizaciones y de la sociedad en general.

²³Management Solutions (2022). AutoML, hacia la automatización de los modelos.

²⁴Management Solutions (2023). Explainable Artificial Intelligence (XAI): desafíos en la interpretabilidad de los modelos.

