



## Nota de prensa

Protección de datos

## Personal investigador de Fisabio diseña un algoritmo que anonimiza la información sensible de los expedientes médicos

- La herramienta permite asegurar que la información con la que trabaja el personal investigador no vulnera los derechos de los pacientes

**València (13.05.2021).** Un proyecto de la Unidad Mixta de Imagen Biomédica de la Fundació per al Foment de la Investigació Sanitària i Biomèdica (Fisabio) y del Centro de Investigación Príncipe Felipe (CIPF) ha desarrollado una herramienta que permite anonimizar textos médicos para que puedan ser usados por investigadores sin vulnerar las leyes de protección de datos.

Los informes clínicos de los pacientes contienen una gran cantidad información útil para los investigadores: pruebas realizadas, medicación del paciente, tiempo de tratamiento, diagnósticos realizados, etc. Irene Pérez-Díez, autora del artículo, recalca que “el tratamiento de este tipo de información ayuda a los investigadores en sus proyectos y contribuye a acelerar el avance científico. Para que los investigadores puedan utilizar esta información, es necesario que el texto esté anonimizado protegiendo la privacidad de los datos personales”.

Este método está basado en categorías (nombre, lugar, número) asociadas a cada unidad de información que luego el algoritmo elimina o cambia por información falsa. Otros métodos anteriores pierden eficacia cuando aparecen errores tipográficos en los textos o expresiones dependientes del contexto. “Nuestro método se basa en el Procesamiento de Lenguaje Natural (NLP), por lo que es sensible al contexto que rodea una palabra determinada; por ejemplo, la IA distingue si la palabra “cabeza” se refiere a una parte del cuerpo o al apellido de una persona”, explica la investigadora.

Los trabajos previos en este campo en castellano no conseguían un nivel perfecto de anonimización, ya que alguna información sensible podía quedar expuesta. Raúl Pérez-Moraga, coautor del artículo, añade que “nuestro método es mucho más robusto y versátil si lo comparamos con métodos basados en reglas fijas o expresiones regulares, ya que estos tienen una fiabilidad deficiente si el informe presenta fallos ortográficos o gramaticales”.

Además, la herramienta se puede trasladar fácilmente a otros idiomas, especialmente si son lenguas derivadas del latín. Según el investigador, “solo haría falta anotar una cantidad suficiente de informes clínicos del idioma en el que se quiera aplicar la metodología. De hecho, el algoritmo ya es capaz de detectar palabras que contiene información sensible tanto en castellano como en valenciano. Esto no ocurre con otros métodos específicos para cada idioma”.

El artículo, titulado ‘De-identifying Spanish medical texts - named entity recognition applied to radiology report’, ha sido publicado en *Journal of Biomedical Sciences* y escrito por Irene Pérez-Díez, Raúl Pérez-Moraga, Adolfo López-Cerdán, Jose-Maria Salinas-Serrano y María de la Iglesia-Vayá, personal investigador de la Unidad Mixta de Imagen de Imagen Biomédica Fisabio-CIPF.

### **Cómo se ha desarrollado el algoritmo**

La metodología de la investigación ha constado de tres fases: anotación, entrenamiento y test. Primero, un equipo de personas expertas revisa los informes clínicos con datos sensibles y anota cada palabra con una etiqueta concreta según de qué tipo sea. “Dividir las palabras a anonimizar en grupos nos permite obtener un mayor rendimiento de los algoritmos de inteligencia artificial”, explica Irene Pérez-Díez.

En una segunda fase de entrenamiento, el equipo experto en Inteligencia Artificial traslada el informe anotado a algoritmos de IA para que estos, “aprendan” los patrones que rodean a las palabras que contiene información sensible. Finalmente se hacen diferentes pruebas para testear el rendimiento. Cuando el algoritmo marca una palabra determinada puede hacer dos cosas, simplemente borrarla o crear información sensible falsa, es decir, si detecta que un nombre de una persona, lo cambiara por otro.

Evaluated los diferentes algoritmos, se elige aquel que haya obtenido un mayor rendimiento para la tarea para la que se ha implementado. Los informes clínicos anonimizados por el algoritmo son también evaluados por equipo de expertos para evitar el filtrado de cualquier información sensible.

Por último, Maria de la Iglesia-Vayá destaca que “el procedimiento desarrollado dentro del marco del proyecto DeepHealth se basa en el principio de protección de datos desde el diseño y por defecto. No sólo se ha desarrollado para operar como un mandato normativo, sino también como una metodología que ayudará en el desarrollo de la transformación digital”.

Esta metodología se presentará en el Hackathon de anonimización que se celebrará a finales de año como parte del Proyecto Europeo DeepHealth.