



## 6101-2. RENDIMIENTO DE LAS EXPRESIONES REGULARES EN EL ANÁLISIS DE INFORMES DE ALTA PRESENTES EN LA HISTORIA CLÍNICA ELECTRÓNICA: EXPRIMIENDO LOS DATOS SECUNDARIOS

Alain García Olea<sup>1</sup>, Ane García Domingo-Aldama<sup>2</sup>, Marcos Merino Prado<sup>2</sup>, Koldo Gojenola Gallettebeitia<sup>2</sup>, Aitziber Atutxa Salazar<sup>2</sup>, Mikel Maeztu Rada<sup>1</sup>, Iván García Díaz<sup>1</sup>, Adrián Costa<sup>1</sup>, Iván Cano<sup>1</sup>, Fernando Díaz<sup>1</sup>, Irene Hernández<sup>1</sup>, Uxue Millet<sup>1</sup>, Ainhoa Etxenike<sup>1</sup> y José Miguel Ormaetxe Merodio<sup>1</sup>

<sup>1</sup>Hospital Universitario de Basurto, Bilbao (Vizcaya), España y <sup>2</sup>Facultad de Ingeniería, Departamento de Lenguajes y Sistemas Informáticos. Universidad del País Vasco, Bilbao (Vizcaya), España.

### Resumen

**Introducción y objetivos:** Este estudio tiene como objetivo investigar cómo la aplicación de expresiones regulares (regex) identifica en los informes de alta de la historia clínica electrónica (HCE) variables no codificadas o mal codificadas en el sistema de informado. Se pretende demostrar cómo esta optimización del proceso de imputación puede mejorar la precisión de los algoritmos predictivos de recurrencia en una población con debut de fibrilación auricular (FA).

**Métodos:** Se realizó un análisis retrospectivo de los informes de alta en formato de texto libre (.txt) de pacientes con debut de FA entre 2015 y 2018. Se entrenó un modelo sobre el que se implementaron expresiones regulares para identificar y extraer información relevante de los informes, especialmente aquella relacionada con las variables codificadas, el debut de FA y su recurrencia. Se compararon los datos obtenidos mediante este método con los datos codificados y se utilizó la comprobación manual de la HCE como *gold standard* comparativo.

**Resultados:** Sobre un *dataset* de 2453 instancias, la aplicación de expresiones regulares sobre los informes de alta resultó en una reducción significativa (58,1%) del número de *missing values* en las variables analíticas codificadas (figura A). Además, gracias a la identificación de valores no codificados como los datos ecocardiográficos (figura B) se observó una mejora sustancial en la integridad y la completud de los datos. La herramienta identificó el debut de FA en el 88,23% de las instancias e identificó en el 61% de los informes datos relativos a la recurrencia o ausencia de recurrencia de FA.



*Porcentaje de valores faltantes en el análisis con regex frente a la codificación hospitalaria (A, %).*  
*Identificación de datos ecocardiográficos en informes de alta (B, %).*

**Conclusiones:** Los resultados de este estudio respaldan la eficacia de utilizar expresiones regulares sobre los informes de alta de la historia clínica electrónica para disminuir los *missing values* en las variables analizadas para estudios a partir de datos secundarios. Además, el análisis de texto libre permite la optimización del

proceso de imputación de datos al identificar variables no codificadas de forma sistemática.