

Resumo do Trabalho de Fin de Grao

(Describa brevemente o traballo a desenvolver, xustificando o interese do mesmo, e indicando obxectivos, descripci3n t3cnica, proceso de desenvolvemento, e medios empregados. Engada tantas li3as como sexa necesario)

La explotaci3n de grandes vol3menes de datos (Big Data) y su uso para la extracci3n de informaci3n y conocimiento o para la implementaci3n de sistemas inteligentes, ha pasado a formar parte de las herramientas b3sicas de traballo de empresas consultoras e instituciones de investigaci3n. La explotaci3n con 3xito de grandes vol3menes de datos supone la necesidad de superar distintas dificultades relacionadas con: *(i)* el volumen, *(ii)* la velocidad en la recepci3n de los datos y *(iii)* la variedad de la informaci3n manejada. Para soportar la variedad de informaci3n manejada *(iii)* se deber3n implementar mecanismos de preprocesamiento que ajusten la informaci3n a un formato 3nico. A mayores, para resolver las problem3ticas del volumen *(i)* y la velocidad de la informaci3n *(ii)*, se debe contar con equipos de gran capacidad computacional que puedan ejecutar las tareas de preprocesamiento de forma que se maximice la cantidad posible de informaci3n por unidad de tiempo. Por tanto, el preprocesamiento de la informaci3n resulta un elemento clave para sacar partido de la explotaci3n de grandes vol3menes de datos, y debe ser realizado de forma que se facilite la definici3n de las operaciones de preprocesamiento que se aplicar3n sobre los datos, se minimicen los errores cometidos en la definici3n del proceso y se maximice la cantidad de informaci3n procesada por unidad de tiempo.

Dadas las caracter3sticas mencionadas en el p3rrafo anterior, este traballo pretende abordar la construcci3n de un framework gen3rico y basado en la definici3n de una tuber3a (pipeline) de traballo de preprocesamiento. El framework ser3 implementado en Java y permitir3 la definici3n de distintas tareas a aplicar sobre los datos as3 como de pipeline de traballo entendido como la secuencia de tareas necesarias para completar el preprocesamiento de datos recibidos de forma heterog3nea. El traballo desarrollado se emplear3 como base para el preprocesamiento de la informaci3n manejada en el proyecto RETOS TIN2017-84658-C2-1-R "Integraci3n de conocimiento sem3ntico para el filtrado de spam basado en contenido" en el cual se pretende la extracci3n de informaci3n de diversas fuentes heterog3neas incluyendo redes sociales, mensajes de correo electr3nico, sitios web diversos, etc.

El desarrollo de la herramienta partir3 del estudio e identificaci3n de los sistemas de preprocesamiento de informaci3n incorporados en herramientas de Big Data y Machine Learning incluyendo herramientas como Mallet (<http://mallet.cs.umass.edu>), TinkerPop (<https://github.com/tinkerpop/pipes>) 3 AIBench (<http://www.aibench.org>). Se intentar3 proponer una soluci3n que permita realizar el preprocesamiento de la informaci3n de una forma sencilla y resolviendo las problem3ticas de las herramientas estudiadas.

El proceso de desarrollo estar3 guiado por la metodolog3a SCRUM.

Para desarrollar el proyecto, se emplear3 el entorno de desarrollo IntelliJ, la herramienta de diagramaci3n Visual Paradigm y el hardware empleado para el desarrollo constar3 de un ordenador port3til DELL XPS15 9560 (i7-7700HQ, 16GB RAM, 512GB SSD, GTX 1050 4GB).