

## Proyecto: Clasificación de Textos en Lenguaje Natural

**Objetivo:** Construir un sistema para la detección de cyber trolls en twitter a partir de un corpus con tweets de cyber trolls y tweets de usuarios normales.

**Aviso importante:** El corpus a tratar contiene expresiones altamente ofensivas. Si algún alumno prefiere trabajar un problema alternativo puede hacerlo comunicándoselo al profesor.

### Contenidos:

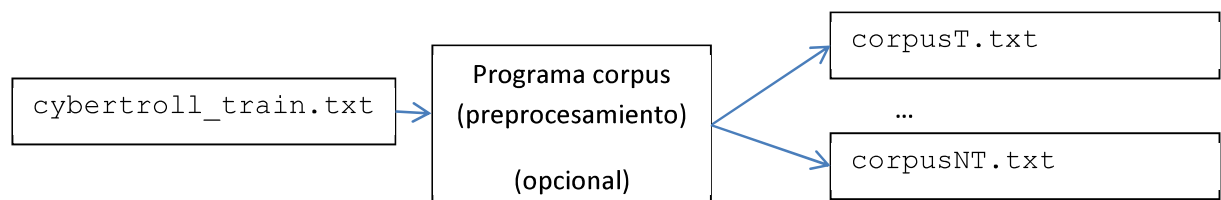
#### *Parte 1 Estimación de probabilidades en el modelo del lenguaje*

En esta parte se estimarán las probabilidades del modelo del lenguaje para las clases troll (T), y notroll (NT).

##### **1.1 Creación de los corpus**

Utiliza el fichero `cybertroll_train.txt` en el campus virtual. Tienes 16435 tweets clasificados en las categorías: troll y not\_troll

Crea 2 corpus con nombre `corpus<T o NT>.txt` con los mensajes de cada categoría. Cada línea del fichero de salida en el corpus debe tener la siguiente estructura:  
<cadena con texto del fichero>



Crea también el fichero `corpustodo.txt` concatenando todos los corpus

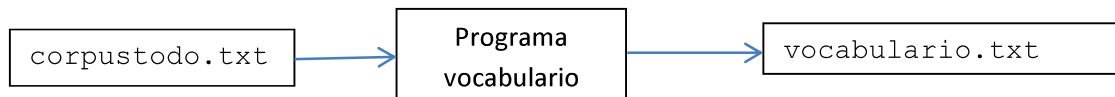
## 1.2 Creación del vocabulario

Halla el vocabulario del problema. Para ello examina el fichero `corpustodo.txt` y obtén las palabras del vocabulario a partir del texto (tokenization).

Debes generar un fichero de salida `vocabulario.txt` con cabecera

Numero de palabras:<Número entero>

Palabra:<cadena>



Las palabras de `vocabulario.txt` estarán ordenadas alfabéticamente.

### Entregable

#### En el Campus Virtual

- Programas:
  - o Corpus (opcional), Vocabulario
- Ficheros:

`corpusT.txt`, `corpusnT.txt`, `corpustodo.txt`, `vocabulario.txt`

### Nota

- Lenguaje de programación libre.