

Práctica con PAN-AP 17.

Asignatura Text Mining en Social Media.

Master Big Data

Elena Tejadillos

elena.tejadillos@gmail.com

Abstract

El objetivo de esta práctica es conseguir un método que partiendo de un tuit, nos permita conocer el género y el país de procedencia del autor del tuit. Para ello se tiene como fuente inicial un conjunto de tuits. De ellos se extraerán una lista de 2000 palabras, las que más frecuentemente aparecen en los tuits, de las cuales se descartarán las stopwords (las palabras que se consideren no relevantes). Estos tuits son procesados y agrupados de manera balanceada en 2 datasets, uno de training y otro de test, en los cuales se estudiará el género y el país de origen de los autores de los tuits. Se aplican 3 modelos diferentes (Support Vector Machine, Naive Bayes con validación cruzada y RandomForest) para ver con cual de todos obtiene un mejor resultado. Finalmente se selecciona el RandomForest como mejor modelo ya que con él se obtiene un *accuracy* más alto, tanto para género como para variedad. Como objetivo final, se deberán superar los umbrales dados. El resultado a superar para reconocer el género es *0.6643* y para reconocer el país es *0.7721*

1. Introducción

Author profiling es el área de estudio que da la posibilidad de conocer rasgos de una persona a partir de los textos que escribe. En el caso concreto de esta práctica, se quiere averiguar el sexo y el país de origen de los autores de un conjunto de tuits. Para realizar el estudio, se ha proporcionado un conjunto de 3200 tuits, que forman parte del dataset PAN-AP'17. Estos tuits se agruparán en dos subconjuntos, uno utilizado para el *training* de los modelos y el otro para el *test*.

Cada dataset se estudiará de manera separada para predecir el género y el país de origen de los autores de los tuits.

2. Dataset

Se proporciona la información necesaria para el dataset PAN-AP'17 a través de un fichero zip. Este fichero contiene los tuits de 300 autores distintos: 200 autores para training y 100 autores para test. Al descomprimirlo, se generan 2 carpetas llamadas *training* y *test*.

En la carpeta *training* se almacenarán 2800 ficheros xml que contienen los tuits a analizar y un fichero *truth.txt* que contiene la lista de títulos de los ficheros xml y el género y país a que corresponden.

La carpeta *test* contiene la misma información pero con tan solo 1400 ficheros xml. Estos tuits son los se van a clasificar dependiendo del género (male/female) o de proceder de uno de los 7 países:

Argentina	Chile
Colombia	España
Mexico	Perú
Venezuela	

Al cargar toda la información en el programa, se tendrán 2 dataset, uno para training y otro para test. Estos datasets se han creado de manera que los datos estén balanceados, lo que significa que habrá el mismo número de tuits para hombres y para mujeres así como para cada uno de los países incluidos en la muestra.

Una vez que ambos datasets están cargados, se hará una limpieza de datos eliminando los signos de puntuación de los textos de los tuits.

4. Resultados experimentales

Tras aplicar los 3 modelos elegidos con una bolsa de 2000 palabras, el mejor resultado obtenido tanto para género como para variedad es el generado por el modelo de Random Forest.

Los resultados son los siguientes:

- Género: 0.7313
- Variedad: 0.8679

Tal y como puede observarse en la siguiente gráfica, los resultados obtenidos por los tres métodos superan, aunque sea levemente, el umbral inicial proporcionado. Pero el Random Forest destaca sobre los tres.

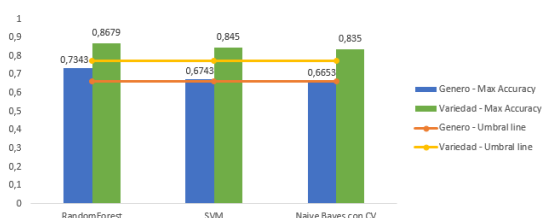


Figura 3: Comparativa de resultados obtenidos y umbrales.

5. Conclusiones y trabajo futuro

Las conclusiones de la práctica han sido que el mejor modelo que podemos utilizar es Random Forest. Pero lo que más ha ayudado a conseguir una mayor mejora del resultado de la predicción, ha sido el aumento de 1000 a 2000 palabras en la bolsa de palabras a utilizar.

Pero este aumento en el número de palabras también ha conllevado a un aumento considerable en el tiempo cálculo de los modelos.

El estudio que se ha hecho de las palabras ha sido muy básico y muy orientado a la variedad, por lo que quedaría pendiente hacer un estudio más en profundidad de las palabras que puedan ser significativas o no por ser utilizadas por los hombres.

También se podría aplicar otros algoritmos como el TF-IDF (*Term frequency – Inverse document frequency*, basada en pesos que permite mejorar la

selección de *stopwords* o como el de S-stemmer, que permite reducir las formas plurales al singular y así calcular la frecuencia de una palabra esté en singular o plural.

References

Mari Vallez y Rafael Pedraza-Jimenez (Universitat Pompeu Fabra) 2007. *El Procesamiento del Lenguaje Natural en la Recuperación de Información Textual y áreas afines*. <https://www.upf.edu/hipertextnet/numero-5/pln.html>