



## Trabajo Práctico Análisis de Texto

Fecha entrega: 29/03/2024

Para la entrega del TP resuelto arme un único archivo (.pdf) y envíelo a través del formulario correspondiente<sup>1</sup>, (el cual se encontrará habilitado hasta la fecha de entrega establecida).

- 1) Escriba un programa que realice análisis léxico sobre la colección **RI-tknz-data**. El programa debe recibir como parámetros el directorio donde se encuentran los documentos y un argumento que indica si se deben eliminar las palabras vacías (y en tal caso, el nombre del archivo que las contiene). Defina, además, una longitud mínima y máxima para los términos. Como salida, el programa debe generar:

- Un archivo (`terminos.txt`) con la lista de términos a indexar (ordenado), su frecuencia en la colección y su DF (*Document Frequency*).

Formato de salida: `$<termino> [ESP] <CF> [ESP] <DF>$`.

Ejemplo:

```
casa 238 3
perro 644 6
...
zorro 12 1
```

- Un segundo archivo (`estadisticas.txt`) con los siguientes datos (un ítem por línea y separados por espacio cuando sean más de un valor):
  - Cantidad de documentos procesados
  - Cantidad de tokens y términos extraídos
  - Promedio de tokens y términos de los documentos
  - Largo promedio de un término
  - Cantidad de tokens y términos del documento más corto y del más largo<sup>2</sup>
  - Cantidad de términos que aparecen sólo 1 vez en la colección
- Un tercer archivo (`frecuencias.txt`), con:
  - La lista de los 10 términos más frecuentes y su CF (Collection Frequency). Un término por línea.
  - La lista de los 10 términos menos frecuentes y su CF. Un término por línea.

<sup>1</sup> Link formulario entrega: <https://forms.gle/NFzYM1cP166rRZkv6>

<sup>2</sup> Mida la longitud del documento en cantidad de tokens.

- 2) Tomando como base el programa anterior, escriba un segundo *tokenizer* que implemente los criterios del artículo de Grefenstette y Tapanainen para definir qué es una “palabra” (o término) y cómo tratar números y signos de puntuación. Además, extraiga en listas separadas utilizando en cada caso una función específica.
- Abreviaturas tal cual están escritas (por ejemplo, Dr., Lic., S.A., etc.)<sup>3</sup>
  - Direcciones de correo electrónico y URLs.
  - Números (por ejemplo, cantidades, teléfonos).
  - Nombres propios (por ejemplo, Villa Carlos Paz, Manuel Belgrano, etc.) y los trate como un único token.

Genere y almacene la misma información que en el caso anterior.

- 3) A partir del programa del ejercicio 1, incluya un proceso de *stemming*<sup>4</sup>. Luego de modificar su programa, corra nuevamente el proceso del ejercicio 1 y analice los cambios en la colección. ¿Qué implica este resultado? Busque ejemplos de pares de términos que tienen la misma raíz pero que el *stemmer* los trató diferente y términos que son diferentes y se los trató igual.
- 4) Sobre la colección CISI<sup>5</sup>, ejecute los *stemmers* Porter y Lancaster provistos en el módulo `nltk.stem`. Compare: cantidad de tokens únicos resultantes, resultado 1 a 1 y tiempo de ejecución para toda la colección. ¿Qué conclusiones puede obtener de la ejecución de uno y otro?
- 5) Escriba un programa que realice la identificación del lenguaje de un texto a partir de un conjunto de entrenamiento<sup>6</sup>. Pruebe dos métodos sencillos:
- Uno basado en la distribución de la frecuencia de las letras.
  - El segundo, basado en calcular la probabilidad de que una letra *x* preceda a otra letra *y* (calcule una matriz de probabilidades con todas las combinaciones).

Compare los resultados contra el módulo Python `langdetect`<sup>7</sup> y la solución provista.

---

<sup>3</sup> Discuta cómo resolvería la extracción de abreviaturas como “NASA”.

<sup>4</sup> Puede usar la librería NLTK (Natural Language Toolkit). Revise qué algoritmos soporta (español e inglés).

<sup>5</sup> [https://ir.dcs.gla.ac.uk/resources/test\\_collections/cisi/cisi.tar.gz](https://ir.dcs.gla.ac.uk/resources/test_collections/cisi/cisi.tar.gz)

<sup>6</sup> Utilice los datos de entrenamiento y de evaluación provistos en: <https://bit.ly/2TqVxIj>

<sup>7</sup> <https://pypi.org/project/langdetect/>



## Propiedades del Texto<sup>8</sup>

- 6) En este ejercicio se propone verificar la predicción de ley de Zipf. Para ello, descargue desde Project Gutenberg el texto del Quijote de Cervantes<sup>9</sup> y escriba un programa que extraiga los términos y calcule sus frecuencias (el programa debe generar la lista ordenada por frecuencia descendente). Calcule la curva de ajuste utilizando la función *Polyfit* del módulo NumPy<sup>10</sup>. Con los datos crudos y los estimados grafique en la notebook ambas distribuciones (haga 2 gráficos, uno en escala lineal y otro en log-log). ¿Cómo se comporta la predicción? ¿Qué conclusiones puede obtener?
- 7) Usando los datos del ejercicio anterior y de acuerdo a la ley de Zipf, calcule la cantidad de palabras que debería haber en el 10%, 20% y 30% del vocabulario. Verifique respecto de los valores reales. Utilice esta aproximación para podar el vocabulario en los mismos porcentajes e indique qué porcentaje de la poda coincide con palabras vacías<sup>11</sup>. Extraiga las palabras podadas que no son *stopwords* y verifique si, a su criterio, pueden ser importantes para la recuperación.
- 8) Codifique un script que reciba como parámetro el nombre de un archivo de texto, tokenize y calcule y escriba a un archivo los pares (`#términos totales procesados`, `#términos únicos`). Verifique en qué medida satisface la ley de Heaps. Grafique en la notebook los ajustes variando los parámetros de la expresión. Puede inicialmente probar con los archivos de los puntos anteriores.

---

<sup>8</sup> Los siguientes ejercicios se deben realizar en una *notebook* IPython para trabajar en un entorno interactivo con capacidades gráficas.

<sup>9</sup> <http://www.gutenberg.org/cache/epub/2000/pg2000.txt> (en UTF-8)

<sup>10</sup> <https://docs.scipy.org/doc/numpy-1.15.0/reference/generated/numpy.polyfit.html>

<sup>11</sup> Use las palabras vacías provistas en la librería NLTK



## **Bibliografía**

- [1] Ricardo Baeza-Yates and Berthier Ribeiro-Neto. Modern Information Retrieval: The Concepts and Technology Behind Search. Addison-Wesley Publishing, USA, 2nd edition, 2008. Cap. 7.
- [2] Gregory Grefenstette and Pasi Tapanainen. What is a word, what is a sentence? problems of tokenization. In Rank Xerox Research Centre, pages 79–87, 1994.
- [3] Le Quan Ha, Darryl Stewart, Philip Hanna, and F Smith. Zipf and type-token rules for the english, spanish, irish and latin languages. Web Journal of Formal, Computational and Cognitive Linguistics, 1 (8):1–12, 1 2006.
- [4] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Sch ¨utze. Introduction to Information Retrieval. Cambridge University Press, New York, NY, USA, 2008. Cap.6.
- [5] Gabriel Tolosa, Fernando Bordignon. Introducci3n a la Recuperaci3n de Informaci3n. Conceptos, modelos y algoritmos b3sicos. Laboratorio de Redes de Datos. UNLu, 2005. Cap. 3.