



Universidad de Granada

PRÁCTICA 1

PARSER DE DOCUMENTOS CON TIKA

Antonio Jiménez Rodríguez

Instrucciones de uso	2
Documentación de código	2
Nube de palabras	3
Ley de Zipf	3

1. Instrucciones de uso

Para facilitar la compilación y ejecución del programa en Java se ha creado un archivo `run.sh` que se hará cargo de estas acciones. Nos colocaremos en el directorio de la práctica y podremos ejecutar:

<code>\$./run.sh -d docs</code>	# Información de los documentos
<code>\$./run.sh -l docs</code>	# Enlaces de cada documento
<code>\$./run.sh -t docs</code>	# Generador de CSV con ocurrencias

2. Documentación de código

Para la realización de esta práctica se han creado dos clases:

→ **Document.java:** La clase *Document* guarda la información sobre el documento que necesitaremos (nombre, tipo, codificación, language y el contenido). La misma clase es la encargada de realizar la extracción de la información haciendo uso de la librería *TIKA*.

Creando un *InputStream* del documento junto con las instancias de la librería *TIKA* necesarias inicializaremos el parser para poder acceder a toda la información. Se usará el *TeeContentHandler* para “unificar” todos los *ContentHandler* que deseamos usar.

Para el tipo del documento accederemos al “Content-Type” del documento gracias a la clase *Metadata*.

`metadata.get("Content-Type")`

Así obtendremos en un *String* el tipo junto a la codificación. Para guardar únicamente el tipo en la propiedad de la clase bastaría con hacer un split de la respuesta.

La codificación del documento, además de poder obtenerla con la clase *Metadata* con el “Content-Type” junto con el tipo podremos acceder a ella con:

`metadata.get("Content-Encoding")`

El cuerpo del documento lo podremos obtener con la instancia de la clase *BodyContentHandler* con el método *toString()*.

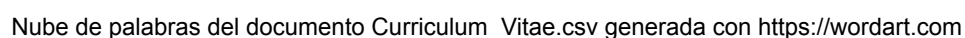
Por último, para el idioma basta con crear una instancia del *LanguageIdentifier* pasándole al constructor el cuerpo del documento, que analizará el texto y devolverá el Lenguaje en el que está escrito.

Los links, al igual que con el cuerpo, hay que llamar al método `getLinks()` de la instancia de la clase `LinkContentHandler`. Este devolverá un `List<Link>` de todos los links del documento.

→ **DocumentsParser.java:** La clase *DocumentsParser* será la encargada de la ejecución. Se han añadido las comprobaciones necesarias para los distintos argumentos que se le deben pasar.

Crearemos el directorio donde guardaremos el CSV y el archivo del cual crearemos el *Writer* donde escribiremos las ocurrencias del documento con el delimitador “;”.

A partir de uno de los documentos CSV generados con la aplicación podemos generar una nube de palabras de forma sencilla.



4. Ley de Zipf

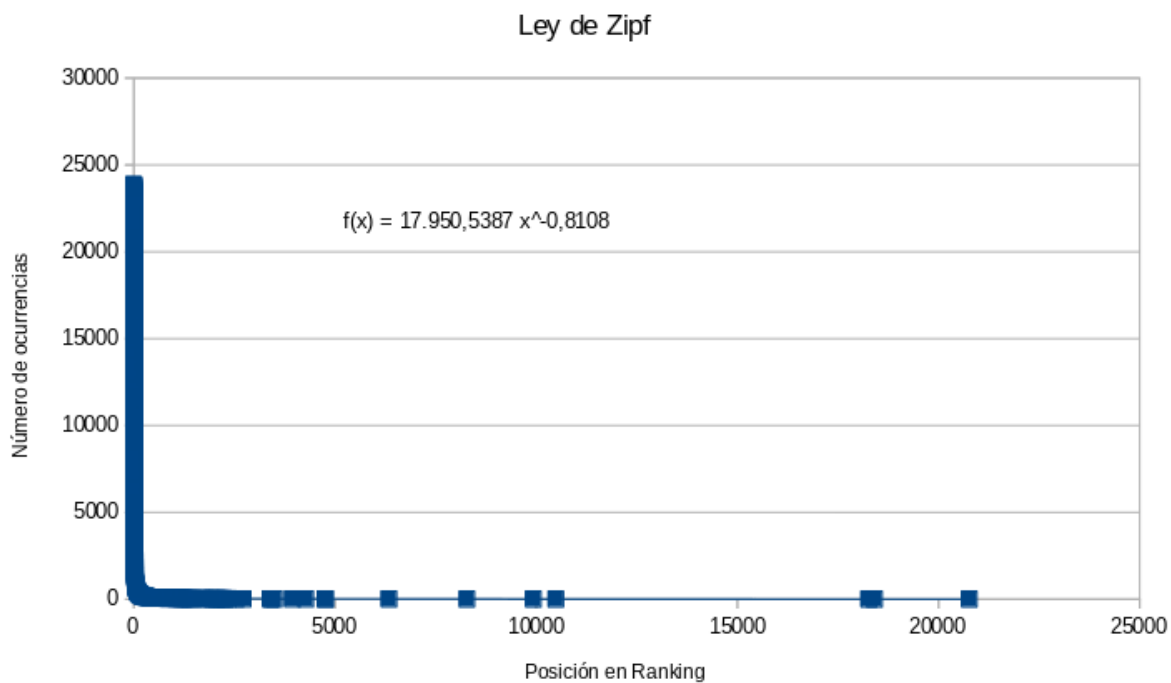


Gráfico realizado con los datos obtenidos del libro *El Quijote*.

De las mediciones del gráfico de dispersión junto con la línea de tendencia podemos recuperar las constantes $k=17959,5387$ y $m=-0,8108$.