

# Práctica 1. Tika

Santiago Álvarez Valdivia

Ioannis Efthymiou

**6 de octubre de 2023**

Recuperación de Información

E.T.S. de Ingenierías Informática y de Telecomunicaciones

**Universidad de Granada**

# Program functionality

-d option -> Table with file name, type, coding and language

```
1 File Summary:
2 -----
3 Filename                               Content-Type                               Encoding                               Language
4 -----
5 file_example_PNG_2500kB.jpg            image/jpeg                                Unknown                               Unknown
6 file-sample_1MB.doc                    application/msword                         Unknown                               ca
7 Grecia.html                           text/html                                 UTF-8                                es
8 Pedro_Paramo-Rulfo_Juan.epub          application/epub+zip                      Unknown                               es
9 WikipediaArticle.html                 text/html                                 UTF-8                                en
10 IvV2y.png                             image/png                                 Unknown                               Unknown
11 file-example_PDF_1MB.pdf               application/pdf                           Unknown                               ca
12 test.txt                              text/plain                                ISO-8859-1                           es
13 file_example_ODS_5000.ods              application/vnd.oasis.opendocument.spreadsheet Unknown                               en
14 file_example_PPT_1MB.pptx              application/vnd.openxmlformats-officedocument.presentationml.presentation Unknown                               ca
15 random.txt                             text/plain                                ISO-8859-1                           en
16 CienAñosDeSoledad.pdf                  application/pdf                           Unknown                               es
```

-l option -> All links that extracted from each document

```
1 File: file_example_PNG_2500kB.jpg
2 -----
3 No links found
4 -----
5 File: file-sample_1MB.doc
6 -----
7 https://products.office.com/en-us/word
8 embedded:image1.emf
9 embedded:image2.jpg
10 -----
11 File: Grecia.html
12 -----
13 /w/load.php?lang=es&modules=codex-search-styles%7Cext.cite.styles%7Cext.tmh.player.styles%7Cext.uls.interlanguage%7Cext.visualE
14 /w/load.php?lang=es&modules=ext.gadget.imagenesinfobox&only=styles&skin=vector-2022
15 /w/load.php?lang=es&modules=site.styles&only=styles&skin=vector-2022
16 /w/load.php?lang=es&modules=noscript&only=styles&skin=vector-2022
17 //upload.wikimedia.org
18 //es.m.wikipedia.org/wiki/Grecia
19 /w/index.php?title=Grecia&action=edit
20 /static/apple-touch/wikipedia.png
21 /static/favicon/wikipedia.ico
```

```

4893 /static/images/footer/poweredby_mediawiki_88x31.png
4894 https://www.mediawiki.org/
4895 -----
4896 File: IvV2y.png
4897 -----
4898 No links found
4899 -----
4900 File: file-example_PDF_1MB.pdf
4901 -----
4902 https://products.office.com/en-us/word
4903 https://products.office.com/en-us/word
4904 https://products.office.com/en-us/word
4905 https://products.office.com/en-us/word
4906 https://products.office.com/en-us/word
4907 https://products.office.com/en-us/word
4908 -----
4909 File: file_example_ODS_5000.ods
4910 -----
4911 No links found
4912 -----
4913 File: file_example_PPT_1MB.pptx
4914 -----
4915 No links found
4916 -----
4917 File: CienAñosDeSoledad.pdf
4918 -----
4919 No links found
4920 -----
4921 -----

```

Note: We used the AutoDetectParser to extract the links for each document, but this does not handle the plain text files, so we did not extract the links from .txt files.

-t option -> count frequency of words in documents, and store them in a csv file per document

Sample output:

```

Counted words for file WikipediaArticle.html. Output in word_count_WikipediaArticle.csv
Counted words for file file-sample_1MB.doc. Output in word_count_file-sample_1MB.csv
Counted words for file IvV2y.png. Output in word_count_IvV2y.csv
Counted words for file file_example_ODS_5000.ods. Output in word_count_file_example_ODS_5000.csv
Counted words for file test.txt. Output in word_count_test.csv
Counted words for file file_example_PNG_2500kB.jpg. Output in word_count_file_example_PNG_2500kB.csv
Counted words for file file-example_PDF_1MB.pdf. Output in word_count_file-example_PDF_1MB.csv
Counted words for file Pedro_Paramo-Rulfo_Juan.epub. Output in word_count_Pedro_Paramo-Rulfo_Juan.csv
Counted words for file Grecia.html. Output in word_count_Grecia.csv
Counted words for file file_example_PPT_1MB.pptx. Output in word_count_file_example_PPT_1MB.csv
Counted words for file CienAñosDeSoledad.pdf. Output in word_count_CienAñosDeSoledad.csv

```

Word art generated from one of the csv files:



# Word Frequency Analysis

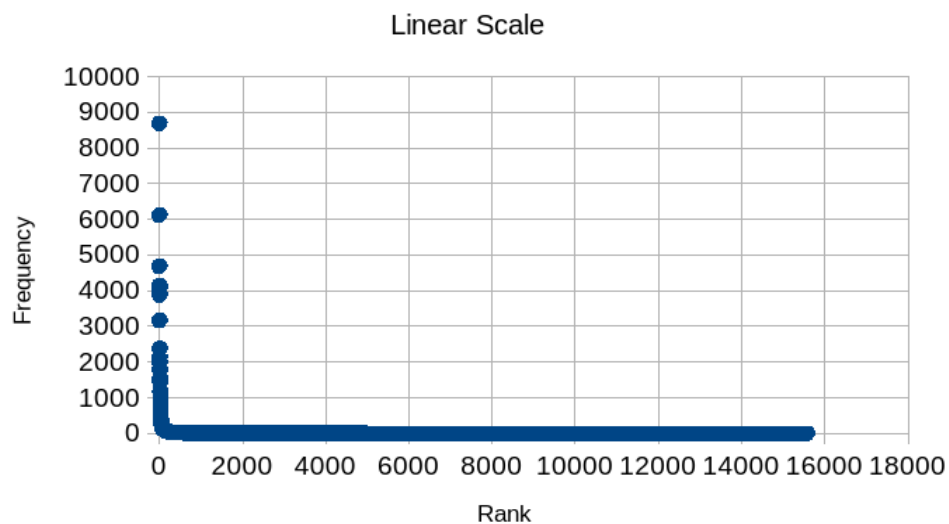
For verifying if a language follows the Zipf's Law, we measured the word frequency for three books, all of them in different languages:

- Cien Años de Soledad, by Gabriel García Márquez, in spanish
- The Collection of Hitchhiker's Guide to the Galaxy books, by Douglas Adams, in British English
- The Illiad, by Homer, in greek

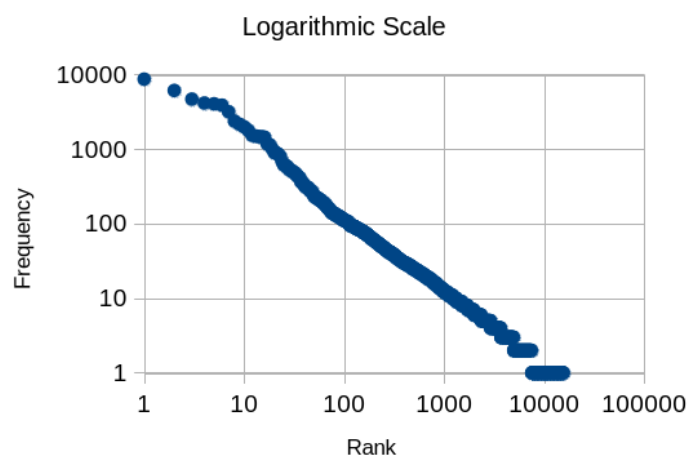
## Spanish

Cien Años de Soledad has 137,865 words, which, when plotted for frequency, look like this.

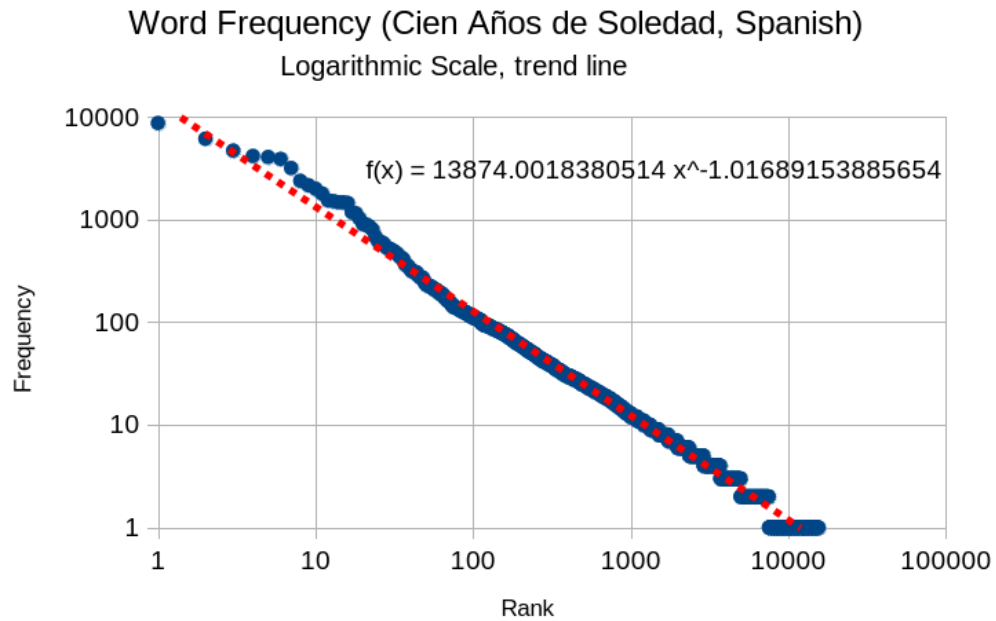
Word Frequency (Cien Años de Soledad, Spanish)



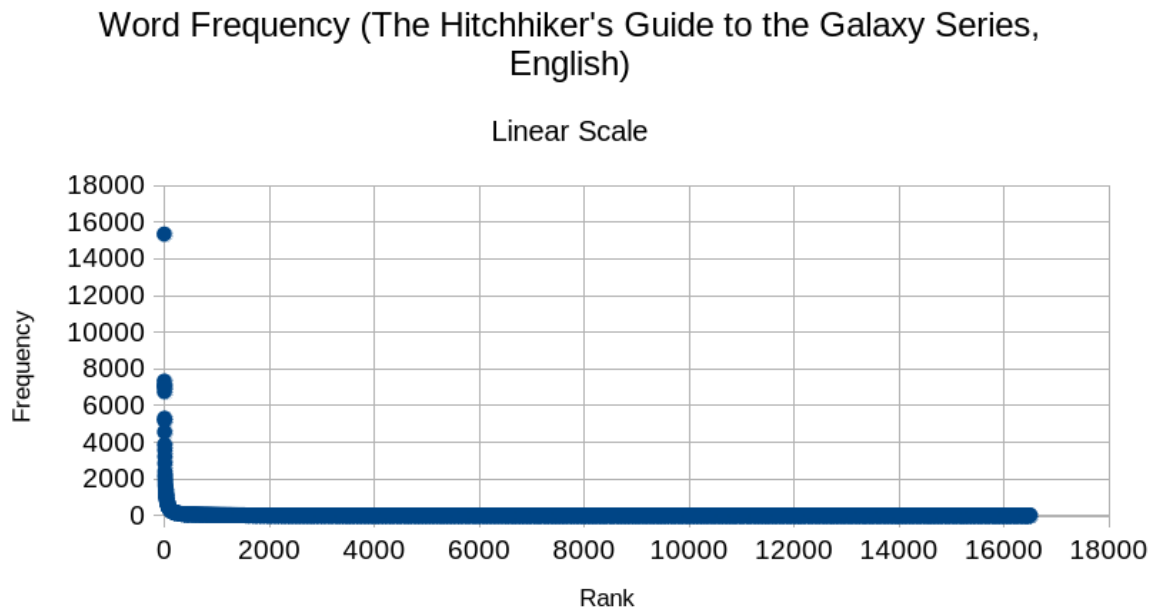
Word Frequency (Cien Años de Soledad, Spanish)



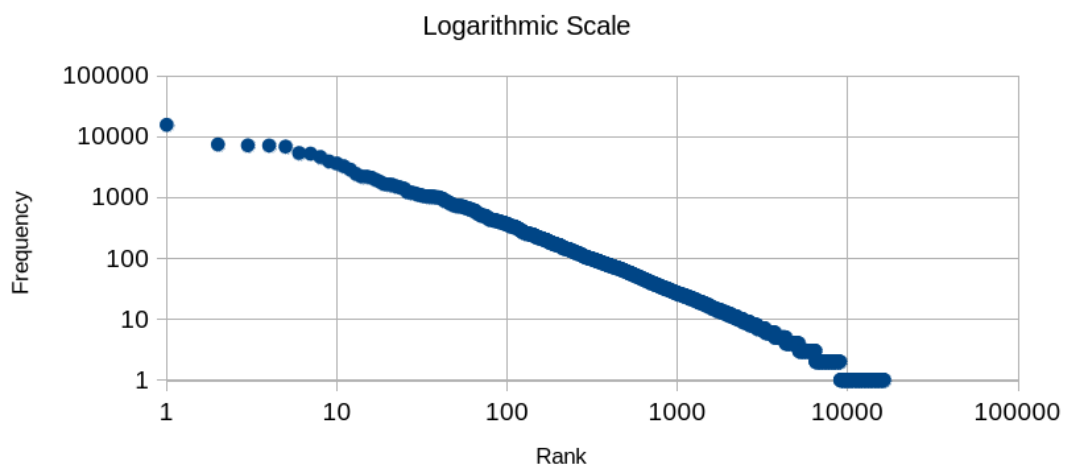
When generating a trend line with a spreadsheet tool, we get a rather close relationship with the actual data. Our language trend is  $F = \frac{k}{R^m}$ , **where k = 13874 and m = 1.01689**



## English

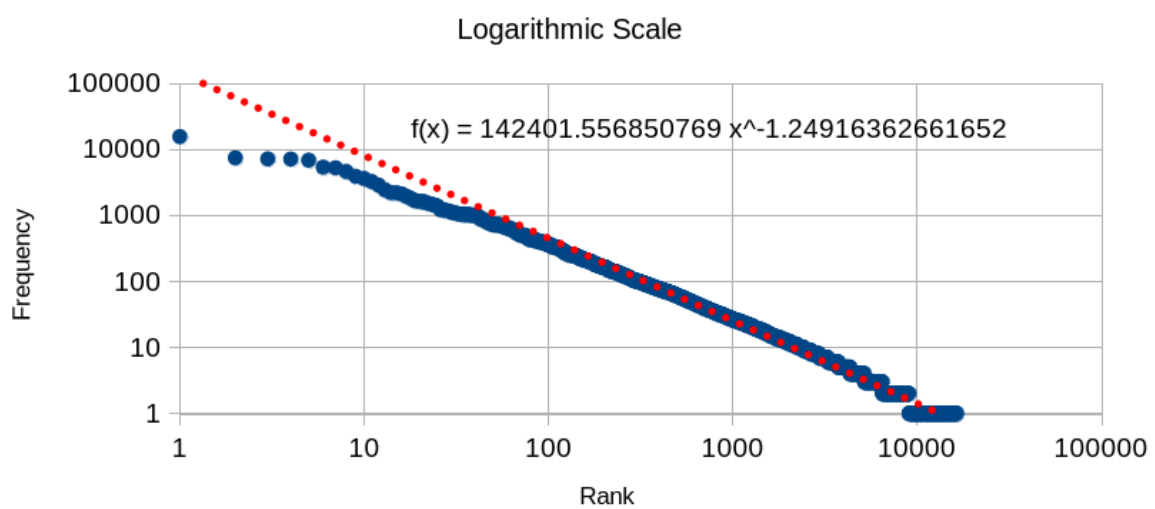


# Word Frequency (The Hitchhiker's Guide to the Galaxy Series, English)



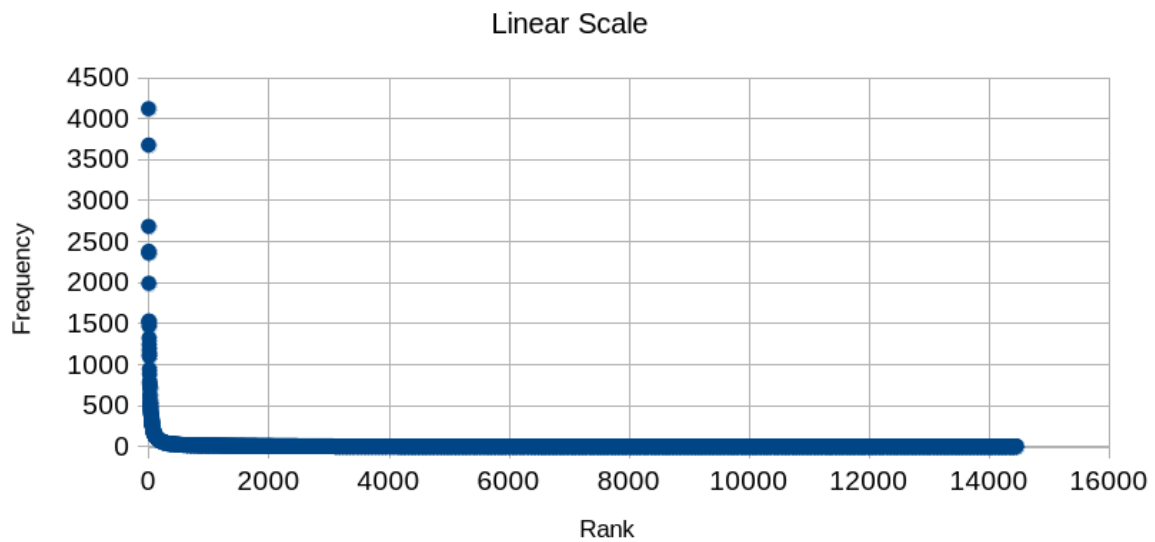
The logarithmic plot trends towards Zipf's Law, but not as closely as our spanish analysis.  
 Our language trend is  $F = \frac{k}{R^m}$ , **where k = 142401 and m = 1.24916**

# Word Frequency (The Hitchhiker's Guide to the Galaxy Series, English)

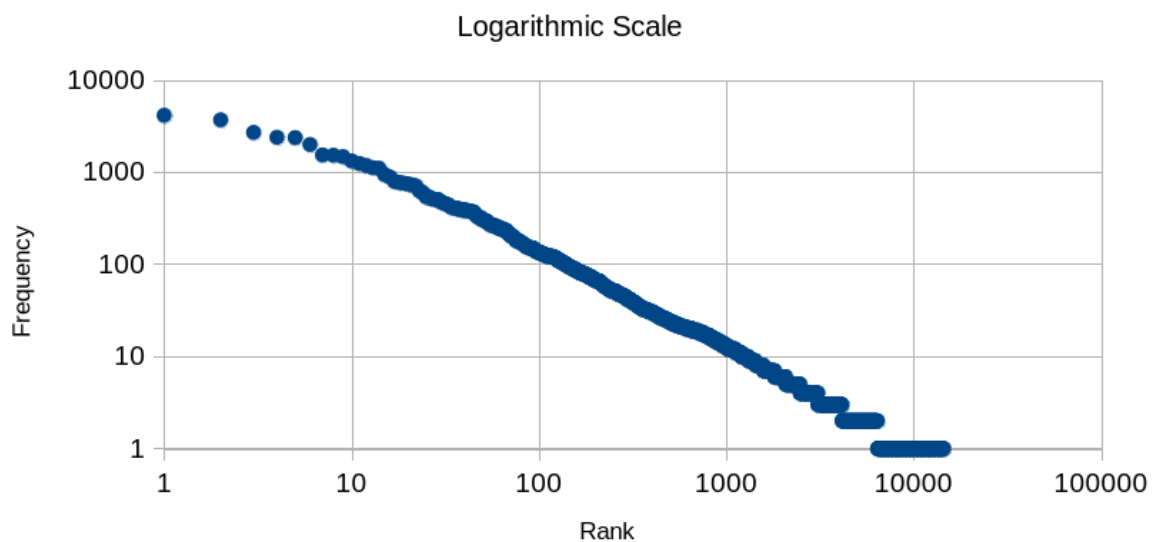


## Greek

Word Frequency (The Illiad, greek)



Word Freequency (The Illiad, greek)



The logarithmic plot trends towards Zipf's Law rather closely. Our language trend is  $F = \frac{k}{R^m}$ , where **k = 14452** and **m = 1.035546**

# Word Freequency (The Illiad, greek)

Logarithmic Scale

