

Zip's

Jordy Joaquín Cuan Robledo ¹, Eduardo Cantorán Flores ¹,
Rafael Pérez Aguirre ¹

¹ Facultad de Ciencias de la Computación, Benemérita Universidad
Autónoma de Puebla, México

El autor de esta fórmula fue George Kingsley Zipf en 1949, él decía que una pequeña cantidad de palabras normalmente tienen una alta frecuencia.

Zip's mide la frecuencia de una palabra que aparece en un texto. Para ello cuenta las repeticiones de una palabra en todo un documento y lo divide entre el total de palabras.

La fórmula es la siguiente:

$$P_n \sim 1/n^a$$

Para mejores resultados se debe emplear un documento suficientemente extenso.

En el lenguaje español las palabras con mayor frecuencia son los artículos y preposiciones.

A continuación se muestra el código:

```
awk '{
1  awk '{
2  gsub(/[/!$%&=?i¿`~""*{,;:()~+[\]\x2E-]/, " ", $0);
3      for (i=1; i<=NF; i++){
4          frecuencia[tolower($i)]++;
5          total++;
6      }
7      next;
8  }
9  END{
10     for (x in frecuencia) {
11         split(x, a, SUBSEP);
12         print a[1], "\t" frecuencia[x], "\t" frecuencia[x]/total;
13     }
14     print "-----"
15     print "Total de palabras = " total;
16 }
17 ' $*
```

Antes de calcular la frecuencia de las palabras contenidas en el documento es necesario limpiar el texto, para ello se utiliza la función gsub la cual quitará los caracteres especiales que delimitamos en cada palabra del texto.

Después recorreremos cada registro insertando cada palabra en una tabla hash aumentando su valor de repetición. Luego para cada palabra de la tabla hash se divide su número de repeticiones entre el total de palabras del texto.

Por último para ejecutarlo lo guardamos en un archivo:

`./zipf.awk textoPrueba.txt > modeloZipf.txt`