

RS - Chapter 04 - Processing Text

Notes

leg ZPT

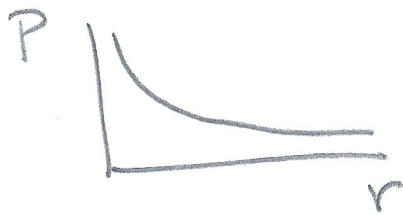
f: frequência = nº ocorrências do b.p. sobre

P: probabilidade = $\frac{f}{N}$
Nº ocorrências totais em
corpus.

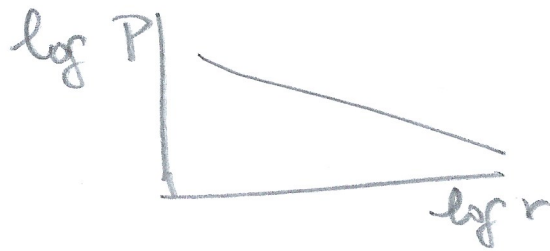
$$f = \frac{N}{r} \quad , \quad \text{equiv.} \quad P = \frac{C}{r}$$

$$P(w) = \frac{C}{r(w)}$$

$$\log P(w) = \log(C) - \log(r(w))$$



$C \approx 0.1$ Inglês



AP89

assistant \rightarrow

$$P = \frac{5.095}{39.749.175} = 0.00013$$

\downarrow
0.013%

$$r \cdot f = k$$

$$\text{rango} \cdot \text{frecuencia} = k$$

$$r_n \cdot n = k$$

$$r_n = k/n$$

$$r_{n+1} \cdot (n+1) = k$$

$$r_{n+1} = \frac{k}{(n+1)}$$

$$\text{ranking última palabra} \cdot 1 = \frac{k}{n(n+1)}$$

$$k = \text{ranking última palabra} = \frac{\text{número total de palabras}}{1}$$

Número de palabras con frecuencia n :

$$r_n - r_{n+1} = \frac{k}{n} - \frac{k}{n(n+1)} = \frac{k}{n(n+1)}$$

$$\frac{\text{Num. palabras frecuencia } n}{\text{Num. total de palabras}} = \frac{1}{k} \frac{k}{n(n+1)} = \frac{1}{n(n+1)}$$

Ex. Ejemplo:
AP89

$$n = 5099 \text{ veces}$$

$$r_n = 1006$$

$$n+1 = 5100 \text{ "}$$

$$r_{n+1} = 1002$$

$$r_n - r_{n+1} = 1006 - 1002 = \underline{4} = \text{nº palabras con frecuencia } 5099$$

Por la predicción

$$\frac{1}{5099 \times 5100}$$

→ población

$$\times 39.799.179$$

$$= (1.53)$$

↓
 $k =$ Número total de

palabras en AP89

Log Heap

$$V = k \cdot n^\beta$$

$$\frac{AP89}{\beta} \quad k = 62.95$$
$$\beta = 0.455$$

$$\frac{Gov2}{\beta} \quad k = 7.34$$
$$\beta = 0.648$$

Estimating Result Set Size

f_a = n^o documentos en los que aparece el término a

Prob. occurrence
término a

$$P(a) = \frac{f_a}{N}$$

Suponendo
independencia
términos

$$f_{abc} = \frac{P(a \cap b \cap c) \cdot N}{N} = P(a) \cdot P(b) \cdot P(c) \cdot N =$$

$$= \frac{f_a}{N} \cdot \frac{f_b}{N} \cdot \frac{f_c}{N} \cdot N = \frac{f_a \cdot f_b \cdot f_c}{N^2}$$

$$P(c | (a \cap b)) = \frac{P(c \cap a \cap b)}{P(a \cap b)}$$

$$P(a \cap b \cap c) = P(a \cap b) \cdot P(c | (a \cap b))$$

a = tropical
c = fish
b = aquarium

Es lo que aproxime por $\frac{P(c|a)}{P(c|b)}$

$$P(a \cap b \cap c) = P(a \cap b) \cdot P(c | b)$$

$$\frac{P(c \cap b)}{P(b)}$$

$$\frac{f_1}{N} = \frac{f_2}{N} \cdot \frac{f_3/N}{f_4/N}$$

$$\uparrow$$

$$P(1) = P(2) \cdot \frac{P(3)}{P(4)} \Rightarrow f_1 = f_2 \cdot \frac{f_3}{f_4}$$

Result Set Estimation

26480 docs processed contains 'aquarium'

De 3000 processed de sus 26480, 258 contains
'tropical' fish 'aquarium'

3000 → 258

26480 → X

$$X = \frac{258 \cdot 26480}{3000} = \frac{C}{S} = 1 \begin{matrix} \text{tropical} \\ \text{fish} \\ \text{aquarium} \end{matrix}$$

C num de docs que contienen las 3 palabras

$$S = \frac{3000}{264800} \rightarrow \text{proporcion del total de docs que tienen palabras}$$

Estimator Collection Site

Suposição de palavras independentes

$$P(a \text{ and } b) = P(a) \cdot P(b)$$

$$\frac{f_{ab}}{N} = \frac{f_a}{N} \cdot \frac{f_b}{N}$$

$$N = \frac{f_a \cdot f_b}{f_{ab}}$$

Estimator Collection nro Google, Bing

lincoln $\rightarrow f_a$

tropical $\rightarrow f_b$

Búsqueda avanzada, f_{ab} lincoln AND tropical

$$N = \frac{f_a \cdot f_b}{f_{ab}}$$

Stemming

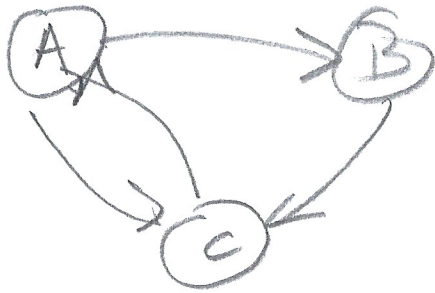
Supply falso negativo: el stemmer no lo detecta, i.e., dice que no es stem y realmente supply es la raíz

up falso positivo: el stemmer si lo detecta pero no es la raíz de UPS

Falsos positivos: ej: organization/organ, el stemmer detecta el stem organ que no es la raíz de organization

Falsos negativos: ej: matrices/matix, el stemmer produce distintos stems pero realmente tienen la misma raíz.

Page Rank



$$PR(C) = \frac{PR(A)}{2} + \frac{PR(B)}{1}$$

$$PR(A) = PR(C)$$

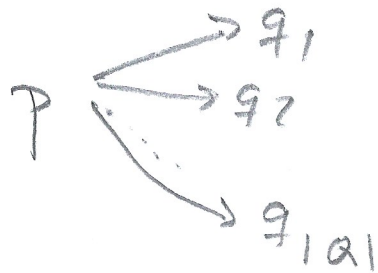
$$PR(B) = \frac{PR(A)}{2}$$

$$\begin{pmatrix} PR(A) & PR(B) & PR(C) \end{pmatrix} \cdot MTP = \begin{pmatrix} PR(A)^T & PR(B)^T & PR(C)^T \end{pmatrix}$$

$$\begin{pmatrix} 0.33 & 0.33 & 0.33 \end{pmatrix} \cdot \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} = \begin{pmatrix} 0.33 & 0.17 & 0.5 \end{pmatrix}$$

...
converge to $\begin{pmatrix} 0.4 & 0.2 & 0.4 \end{pmatrix}$

l.11: $\pi_i \leftarrow \frac{\lambda}{|P|}$ probabilidad de cada página por teleputing



l.17 $R_q = R_q + (1-\lambda) \frac{I_P}{|Q|}$ \rightarrow propaga

el PR de P, I_P , a los págs q a los que apunta P de forma equiprobable y a la cantidad total ponderada por $(1-\lambda)$

l.21. Si P no tiene páginas salientes, $|Q|=0$, su PageRank se propaga de forma equiprobable a todas las páginas del sitio.