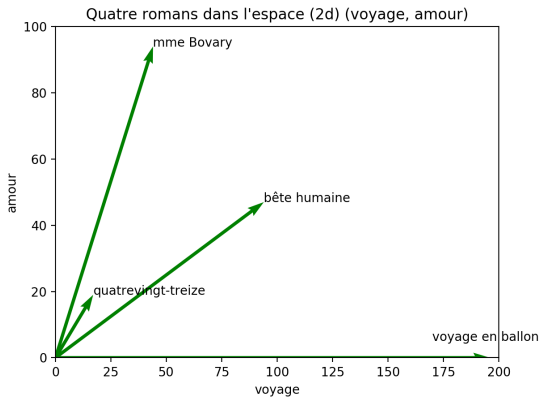


- *If A and B have almost identical environments we say that they are synonyms. (Harris, 1954)*
- *You shall know a word by the company it keeps. (Firth, 1957)*
- *The degree of semantic similarity between two linguistic expressions A and B is a function of the similarity of the linguistic contexts in which A and B can appear. (Lenci, 2008)*

	QuatreVT 119 Kw	Voyage Bal 82 kw	Bête Hum. 128 kw	Mme Bovary 117 kw
bataille	35	4	6	2
clair	105	26	96	52
facile	12	19	6	10
politique	11	0	9	5
voyage	17	196	94	44
idiot	2	1	2	6
amour	19	0	47	94

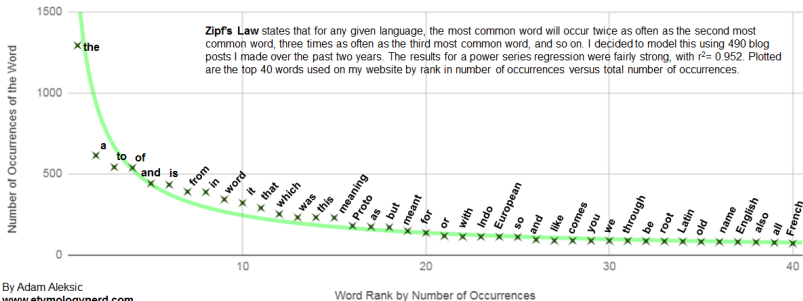
**Table** – Matrice terme-documents pour quelques mots et 4 romans (Quatrevingt-treize (Hugo); Le voyage en ballon (Verne); La bête humaine (Zola); Mme Bovary (Flaubert)).

	QuatreVT	Voyage Bal	Bête Hum.	Mme Bovary
voyage	17	196	94	44
amour	19	0	47	94



# Loi de Zipf

## Zipf's Law and My Blog on Language



By Adam Aleksic  
[www.etymologynerd.com](http://www.etymologynerd.com)

## tf-idf : formules

$$tf(t, d) = \begin{cases} \log(1 + f_{t,d}) & \text{si } f_{t,d} > 0 \\ 0 & \text{sinon} \end{cases}$$

$$idf(t, D) = \log\left(\frac{|D|}{df(t)}\right)$$

$$td-idf(t, d) = tf(t, d, D) \times idf(t, D)$$

bataille (35,2)  
politique (11,5)  
amour (19,94)  
voyage (17,44)

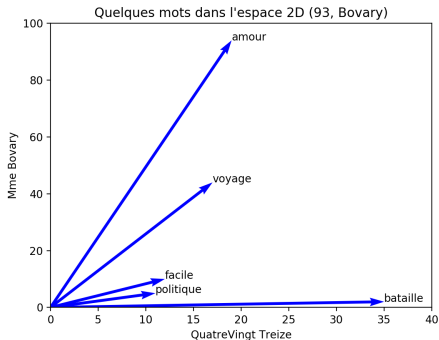


Figure – représentation graphique des mots dans le plan (93, Bovary)



	bataille	voyage	homme	femme
arriver	246	470	1 819	890
tomber	100	83	1 205	384
habiller	2	4	339	384
mourir	180	116	339	1 088
	55 331	208 520	668 289	346 093

Table – Matrice Terme-Terme obtenue dans frWaC. Contextes en ligne.



	arriver	tomber	habiller	mourir	
bataille	246	100	2	180	55 331
voyage	470	83	4	116	208 520
homme	1 819	1 205	339	1 499	668 289
femme	890	660	384	1 088	346 093

**Table** – Matrice Terme-Terme obtenue dans frWaC. Contextes en colonnes.

$f$	aardvark	computer	data	pinch	result	sugar
apricot	0	0	0	1	0	1
pineapple	0	0	0	1	0	1
digital	0	2	1	0	1	0
information	0	1	6	0	4	0

$f$	computer	data	pinch	result	sugar	$\Sigma$
apricot	0	0	1	0	1	2
pineapple	0	0	1	0	1	2
digital	2	1	0	1	0	4
information	1	6	0	4	0	11
$\Sigma$	3	7	2	5	2	19

$p_{ij}$	computer	data	pinch	result	sugar	$p_{a*}$
apricot	0/19	0/19	1/19	0/19	1/19	2/19
pineapple	0/19	0/19	1/19	0/19	1/19	2/19
digital	2/19	1/19	0/19	1/19	0/19	4/19
information	1/19	6/19	0/19	4/19	0/19	11/19
$p_{*b}$	3/19	7/19	2/19	5/19	2/19	19/19

$p_{ij}$	computer	data	pinch	result	sugar	$p_{a*}$
apricot	0,00	0,00	0,05	0,00	0,05	0,11
pineapple	0,00	0,00	0,05	0,00	0,05	0,11
digital	0,11	0,05	0,00	0,05	0,00	0,21
information	0,05	0,32	0,00	0,21	0,00	0,58
$p_{*b}$	0,16	0,37	0,11	0,26	0,11	

<i>ppmi</i>	computer	data	pinch	result	sugar
apricot	-	-	2,25	-	2,25
pineapple	-	-	2,25	-	2,25
digital	1,66	0,00	-	0,00	-
information	0,00	0,57	-	0,47	-

$f$	computer	data	pinch	result	sugar
apricot	2	2	3	2	3
pineapple	2	2	3	2	3
digital	4	3	2	3	2
information	3	8	2	6	2

$p_{ij}$	computer	data	pinch	result	sugar	$p_{a*}$
apricot	0,03	0,03	0,05	0,03	0,05	0,20
pineapple	0,03	0,03	0,05	0,03	0,05	0,20
digital	0,07	0,05	0,03	0,05	0,03	0,24
information	0,05	0,14	0,03	0,10	0,03	0,36
$p_{*b}$	0,19	0,25	0,17	0,22	0,17	



<i>ppmi</i>	computer	data	pinch	result	sugar
apricot	-	-	2,25	-	2,25
pineapple	-	-	2,25	-	2,25
digital	1,66	0,00	-	0,00	-
information	0,00	0,57	-	0,47	-
<i>ppmi</i> [add2]	computer	data	pinch	result	sugar
apricot	0,00	0,00	0,56	0,00	0,56
pineapple	0,00	0,00	0,56	0,00	0,56
digital	0,62	0,00	0,00	0,00	0,00
information	0,00	0,58	0,00	0,37	0,00

Figure – Comparaison de mesures de *ppmi* sur l'exemple de (Jurafsky & Martin, 2019), avec ou sans lissage laplacien +2

## Références

- Firth, John Rupert. 1957. *Papers in Linguistics (1934-1951)*. Oxford : Oxford University Press.
- Harris, Zellig S. 1954. Distributional structure. *Word*, 10(2-3), 146–162.
- Jurafsky, Daniel, & Martin, James H. 2019. *Speech and Language Processing : An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics*. Prentice-Hall. drafts of August 29, 2019. Chap. 6 Vector Semantics.
- Lenci, Alessandro. 2008. Distributional semantics in linguistic and cognitive research. *Italian journal of linguistics*, 20(1), 1–31.