

01

Ziff's law +

It states that frequency of words at rank α is inversely proportional to rank of word in corpus.

$$f(\alpha) = \frac{f(1)}{\alpha}$$

$f_{(1)}$ = Frequency of most frequent words (rank₁)

i) If rank of word is doubled. what should be expected frequency of word.

$\gamma_2 = 2\gamma_1$ (6.128. 58) $\gamma_1 = 1$ mmol, word 22

$$f(\sigma_1) = \frac{f(\sigma)}{\sigma_1} \text{ (approx.)} \quad 1 + [0.0] - i0b = (1, i0b)$$

$$f(x_2) = f(1)$$

$$\underline{f(1)} \Rightarrow f(1) \geq \underline{500} = 250$$

9 142578 2 2 2 0 TAR & S

o utraque to 2 o utraque 2

(i) $N = 1000$ words; $R(t) = 50$; $f(t) = ?$

9902-0216 2000 Aug 22

f(x) graph for UR | 0 0 0 . accuracy

8- 615 O MAGUIRA & C.

Digitized by srujanika@gmail.com

$$f(10) = \frac{50}{10} = 5$$

$$f(10) = \frac{50}{10} = 5$$

2) Min edit Distance

	S	A	T	U	R	D	A	Y
S	0	1	2	3	4	5	6	7
U	1	1	2	3	4	5	6	7
N	2	2	2	3	4	5	6	7
D	3	3	3	3	4	3	4	5
A	4	3	4	5	5	4	3	4
Y	5	4	4	6	6	5	4	3

We know formula: $d[i-1] = T[j-1]$

$$d[i, j] = \begin{cases} d[i-1][j-1] + 1 & \text{(Delete) [D]} \\ d[i][j-1] + 1 & \text{(Insert) [I]} \\ d[i-1][j-1] & \text{Substitute [S]} \end{cases}$$

Seq	D	I	S		D	I	S
i) $S \rightarrow S$	0	0	0	ii) $SU \rightarrow S$	0	1	P
$S \rightarrow SA$	0	1	0	$SU \rightarrow SA$	0	0	I
$S \rightarrow SAT$	0	2	0	$SU \rightarrow SAT$	0	1	I
$S \rightarrow SATU$	0	3	0	$SU \rightarrow SATU$	0	2	I
$S \rightarrow SATUR$	0	4	0	$SU \rightarrow SATUR$	0	3	I
$S \rightarrow SATURD$	0	5	0	$SU \rightarrow SATURD$	0	4	I
$S \rightarrow SATURDA$	0	6	0	$SU \rightarrow SATURDA$	0	5	I
$S \rightarrow SATURDAY$	0	7	0				

	D	I	S
iii) SUN \rightarrow S	2	0	0

SUN \rightarrow SA	1	0	1
----------------------	---	---	---

SUN \rightarrow SAT	0	0	2
-----------------------	---	---	---

SUN \rightarrow SATU	0	1	2
------------------------	---	---	---

SUN \rightarrow SATUR	0	2	2
-------------------------	---	---	---

SUN \rightarrow SATURD	0	3	2
--------------------------	---	---	---

SUN \rightarrow SATURDA	0	4	2
---------------------------	---	---	---

D	I	S
---	---	---

IV) SUND \rightarrow S	0	3	0
--------------------------	---	---	---

SUND \rightarrow SA	0	2	0
-----------------------	---	---	---

SUND \rightarrow SAT	0	0	2
------------------------	---	---	---

SUND \rightarrow SATU	0	2	0
-------------------------	---	---	---

SUND \rightarrow SATUR	0	2	3
--------------------------	---	---	---

SUND \rightarrow SATURD	0	2	0
---------------------------	---	---	---

SUND \rightarrow SATURDA	0	3	1
----------------------------	---	---	---

SUND \rightarrow SATURDAY	0	4	1
-----------------------------	---	---	---

D	I	S
---	---	---

V) SUND A \rightarrow S	4	5	0
---------------------------	---	---	---

SUND A \rightarrow SA	3	0	0
-------------------------	---	---	---

SUND A \rightarrow SAT	3	1	0
--------------------------	---	---	---

SUND A \rightarrow SATU	3	2	0
---------------------------	---	---	---

SUND A \rightarrow SATUR	2	2	1
----------------------------	---	---	---

SUND A \rightarrow SATURD	1	2	1
-----------------------------	---	---	---

SUND A \rightarrow SATURDA	0	2	1
------------------------------	---	---	---

SUND A \rightarrow SATURDAY	0	3	1
-------------------------------	---	---	---

D	I	S
---	---	---

VI) SUNDAY \rightarrow S	5	0	0
----------------------------	---	---	---

SUNDAY \rightarrow SA	4	0	0
-------------------------	---	---	---

SUNDAY \rightarrow SAT	3	0	1
--------------------------	---	---	---

SUNDAY \rightarrow SATU	4	2	0
---------------------------	---	---	---

SUNDAY \rightarrow SATUR	3	2	1
----------------------------	---	---	---

SUNDAY \rightarrow SATURD	2	2	1
-----------------------------	---	---	---

SUNDAY \rightarrow SATURDA	1	2	1
------------------------------	---	---	---

SUNDAY \rightarrow SATURDAY	0	2	1
-------------------------------	---	---	---

Q.3

Given Sentence :-

S₁ : "The cat sat on the mat"

S₂ : "The dog jumped over the fence."

S₃ : The cat and the dog played together"

Vocab

Tokenize :-

S₁ : "the", "Cat", "sat", "on", "the", "mat"

S₂ : "the", "dog", "jumped", "over", "the", "fence"

S₃ : "The", "Cat", "and", "the", "dog", "played",
"together"

BOW :-

Vocab : [and, the, jumped, mat, on, played, cat, dog,
together, *sat, over, fence]

words

doc1

doc2

+ doc3

and

the

jumped

mat

On

played

cat

dog

together

sat

BOW : (Bag of words)

FACT-IT (ii)

Word	Doc 1	Doc 2	Doc 3
and	0	0	1
cat	1	0	1
dog	1	1	1
fence	1	1	0
jumped	0	1	0
mat	1	0	0
on	1	0	0
over	0	1	0
played	0	0	1
sat	1	0	0
the	2	1	2
together	0	0	1

Doc 1 : [0 1 0 0 0 1 0 1 0 0 1 2 0]

Doc 2 : [0 0 1 1 1 0 1 0 1 0 0 2 0]

Doc 3 : [1 1 1 0 0 0 0 0 1 0 0 2 1]

Now,

ii) TF-IDF

we know that;

TF → Term frequency

$TF = \frac{\text{No of times term } t \text{ in document } D}{\text{Total no of terms in Document } D}$

word	D ₁	D ₂	D ₃	
and	0	1/0	1/7	
eat	1/6	0	1/7	
dog	0	1/6	1/7	
fence	0	1/6	0	
jumped	0	1/6	0	
mat	1/6	0	0	
on	1/6	1/0	0	
over	0	1/6	0	
played	0	0	1/7	
sat	1/6	0	0	
the	2/6	2/6	2/7	
together	0	0	1/7	

IDF

$$IDF = \lg \left(\frac{N}{1 + DF(t)} \right) = \frac{\text{Total no of documents}}{\text{No of documents containing term t}}$$

Document Frequency (DF)

Word	DFG	IDF	IDF	TF-IDF		
				D ₁	D ₂	D ₃
and	1	$\lg(3/2)$	0.405	0	1	0.057
cat	2	$\lg(3/3)$	0	0	0	0
dog	2	$\lg(3/3)$	0	0	0	0
fence	1	$\lg(3/2)$	0.405	0	0.067	0
jumped	1	$\lg(3/2)$	0.405	0	0.067	0
mat	1	$\lg(3/2)$	0.405	0.067	0	0
on	1	$\lg(3/2)$	0.405	0.067	0	0
over	1	$\lg(3/2)$	0.405	0	0.067	0
played	1	$\lg(3/2)$	0.405	0	0	0.057
sat	1	$\lg(3/2)$	0.405	0.067	0	0
the	3	$\lg(3/3)$	-0.287	-0.095	-0.095	-0.082
together	1	$\lg(3/2)$	0.405	0	0	0.057

$$D_1 = [0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0]$$

$$D_2 = [0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$$

$$D_3 = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 1]$$

Similarity is 0 for all combinations

$$D_1 \& D_2 = 0$$

$$D_1 \& D_3 = 0$$

$$D_2 \& D_3 = 0$$

Similarity of Doc :-

Doc 1 & Doc 2

$$\text{Magnitude of Doc 1} : \sqrt{1+1+1+1+4} = \sqrt{8}$$

$$\text{Magnitude of Doc 2} : \sqrt{1+1+1+1+4} = \sqrt{8}$$

$$\text{Magnitude of Doc 3} : \sqrt{1+1+1+1+4} = \sqrt{9} = 3$$

Doc 1 & Doc 2

$$\text{Sim 1} : \frac{4}{\sqrt{8}\sqrt{8}} = \frac{4}{8} = 0.5 \approx 0.5$$

Doc 1 & Doc 3

$$\text{Sim 2} : \frac{(1+4)}{\sqrt{8}\sqrt{9}} = \frac{5}{\sqrt{8}\sqrt{9}} = 0.589 \approx 0.6$$

Doc 2 & Doc 3

$$\text{Sim 3} : \frac{(1+4)}{\sqrt{8}\sqrt{9}} = \frac{5}{\sqrt{8}\sqrt{9}} = 0.589 \approx 0.6$$

Q.3

Doc 1: "The quick brown fox jumps over the lazy dog" 9

Doc 2: "The lazy cat sleeps on the sunny mat". 8

Doc 3: "The quick brown fox likes to jump high". 8

(without removing stop words)

Tokenerized:

Doc 1: ['the', 'quick', 'brown', 'fox', 'jumps', 'over', 'the', 'lazy', 'dog']

Doc 2: ['the', 'lazy', 'cat', 'sleeps', 'on', 'the', 'sunny', 'mat']

Doc 3: ['the', 'quick', 'brown', 'fox', 'likes', 'to', 'jump', 'high']

Vocab:

Un O: [the, quick, brown, fox, jumps, over, lazy, dog, cat, sleeps, on, sunny, mat, likes, to, jump, high]

O: [brown, cat, dog, fox, high, jump] jumps, lazy, likes, mat, on, over, quick, sleeps, sunny, the, to]

Bag of words vector:

D₁: [1 0 0 1 0 1 1 0 0 0 1 1 0 0 2 0]

D₂: [0 1 0 0 0 0 1 0 1 1 0 0 1 1 2 0]

D₃: [1 0 0 1 1 1 0 1 0 0 0 1 0 0 1 1]

$$\text{Mag of } D_1 : \sqrt{1+1+1+1+1+1+1+4} = \sqrt{11}$$

$$\text{Mag of } D_2 : \sqrt{1+1+1+1+1+1+9} = \sqrt{10}$$

$$\text{Mag of } D_3 : \sqrt{1+1+1+1+1+1+1+1} = \sqrt{8}$$

Doc 1 & Doc 2:

$$\text{Sim 1} = \frac{5}{\sqrt{11}\sqrt{10}} = 0.476$$

Doc 1 & Doc 3:

$$\text{Sim 2} = \frac{6}{\sqrt{11}\sqrt{8}} = 0.639$$

Doc 2 & Doc 3:

$$\text{Sim 3} = \frac{2}{\sqrt{10}\sqrt{8}} = 0.2236$$

TF-IDF:

Words	TF			DF	IDF	IDF	TF-IDF		
	D ₁	D ₂	D ₃				D ₁	D ₂	D ₃
brown	1/4	0	1/8	2	lg(3/3)	0	0	0	0
cat	0	1/8	0	1	lg(3/2)	0.176	0	0.022	0
dog	1/4	0	0	1	lg(3/2)	0.176	0.019	0	0
for	1/6	0	1/8	2	lg(3/3)	0	0	0	0
high	0	0	1/8	1	lg(3/2)	0.176	0	0	0.022
jump/s	1/4	0	1/8	2	lg(3/3)	0	0	0	0
lazy	1/6	1/8	0	2	lg(2/3)	0	0	0	0
likes	0	0	1/8	1	lg(3/2)	0.176	0	0	0.022
mat	0	1/8	0	1	lg(3/2)	0.176	0	0.022	0
on	0	1/8	0	1	lg(3/2)	0.176	0	0.022	0
over	1/6	0	0	1	lg(3/2)	0.176	0.019	0	0
quick	1/6	0	1/8	2	lg(3/3)	0	0	0	0
sleeps	0	1/8	0	1	lg(3/2)	0.176	0	0.022	0
sunny	0	1/8	0	1	lg(3/2)	0.176	0	0.022	0
the	2/6	2/8	1/8	3	lg(3/5)	-0.124	-0.027	-0.031	-0.015
to	0	0	1/8	1	lg(3/2)	0.176	0	0	0.022

Similarity,

D₁ & D₂

$$\text{Sim 1: } \frac{0.00041}{(0.038)(0.058)} = 0.186$$

$$D_1 > D_2 = .$$

$$|D_1| = 0.038$$

$$|D_2| = 0.058$$

$$|D_3| = 0.0409$$

D₁ & D₃

$$\text{Sim 2: } \frac{0.00041}{(0.038)(0.0409)} = 0.263$$

D₂ & D₃

$$\text{Sim 3: } \frac{0.00047}{(0.058)(0.0409)} = 0.198$$