



# IIC 3800 Tópicos en CC NLP

<https://github.com/marcelomendoza/IIC3800>

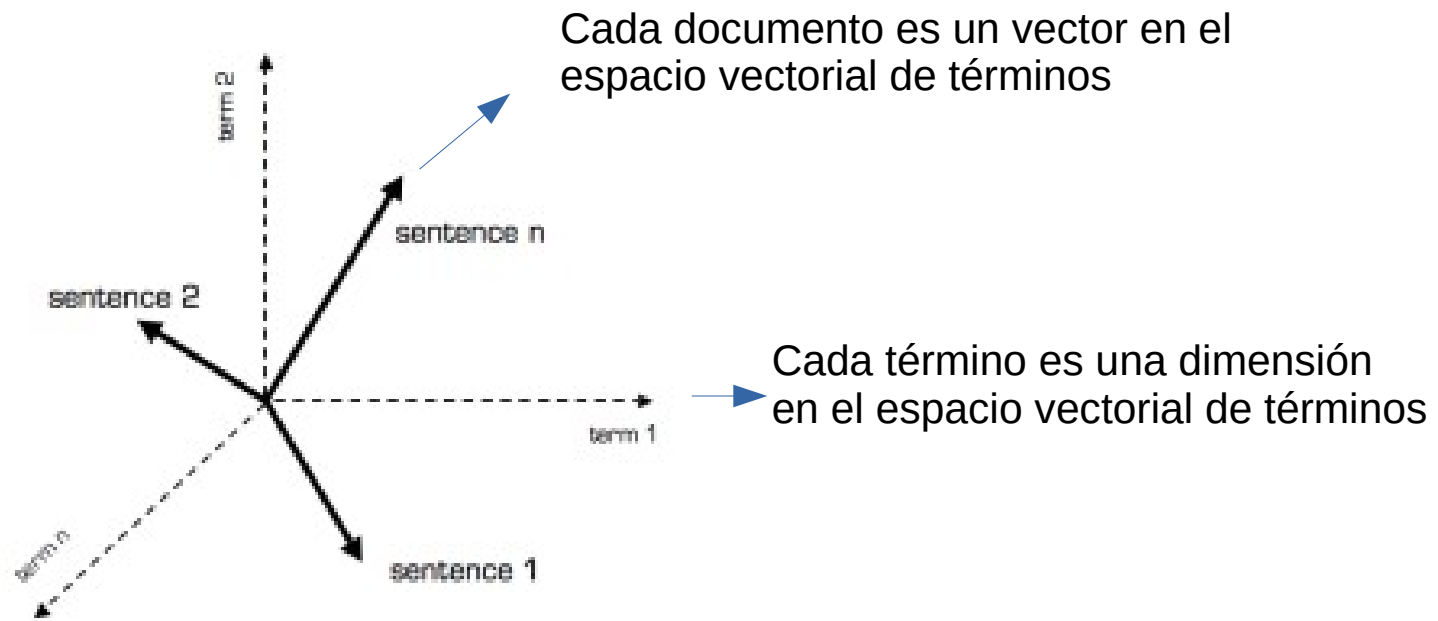
# - VECTORIZACIÓN DE DOCUMENTOS Y RANKING -

## Matriz términos-documentos

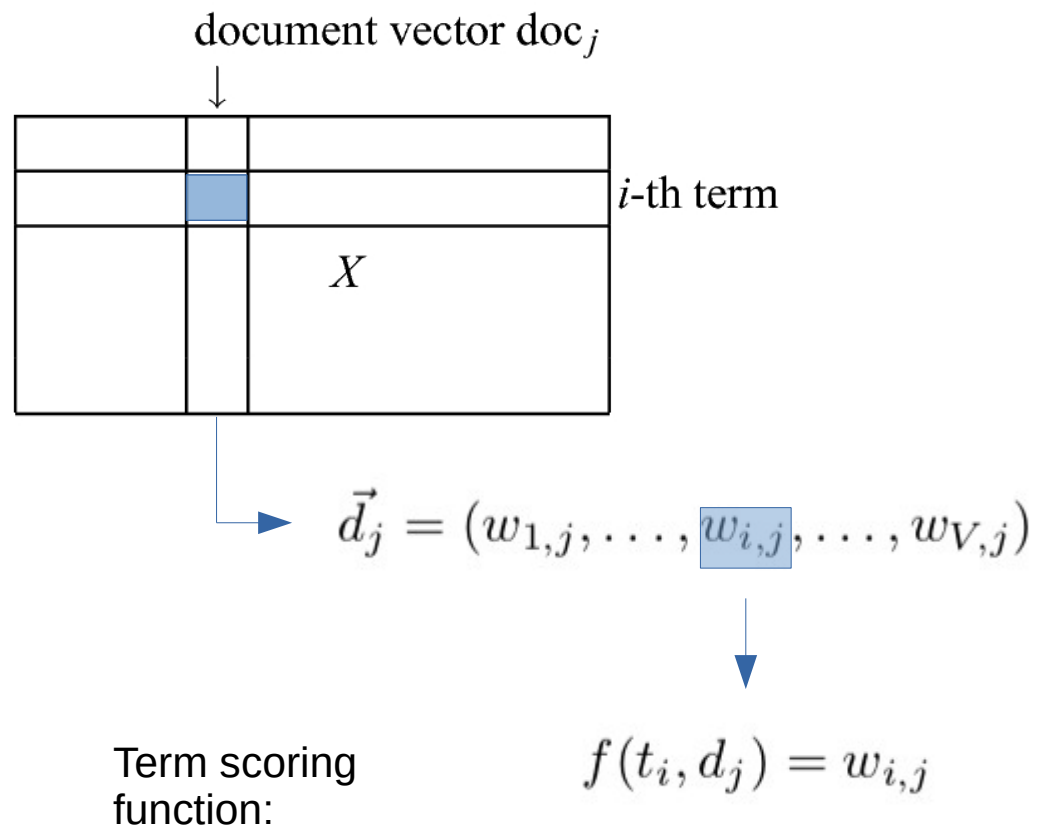
documentos	Antonio y Cleopatra	Julio Cesar	La Tempestad	Hamlet	Otelo	Macbeth	...
Antonio	157	73	0	0	0	1	
Brutus	4	157	0	2	0	0	
Cesar	232	227	0	2	1	0	
Calpurnia	0	10	0	0	0	0	
Cleopatra	57	0	0	0	0	0	
...		...					
términos							

vectorización

## Vector-space model



## Vector-space model



BM25

$f_{i,j}$  : # occs. de ti en dj

$N$  : # docs

$n_i$  : # docs donde ti ocurre

$l(d_j)$  : # tokens en dj

$l_{avg}$  : largo promedio

$$w_{i,j} = \frac{f_{i,j} \cdot (k_1 + 1)}{k_1 \cdot \left[ (1 - b) + b \cdot \frac{l(d_j)}{l_{avg}} \right] + f_{i,j}} \cdot \log \left( \frac{N}{n_i} \right)$$

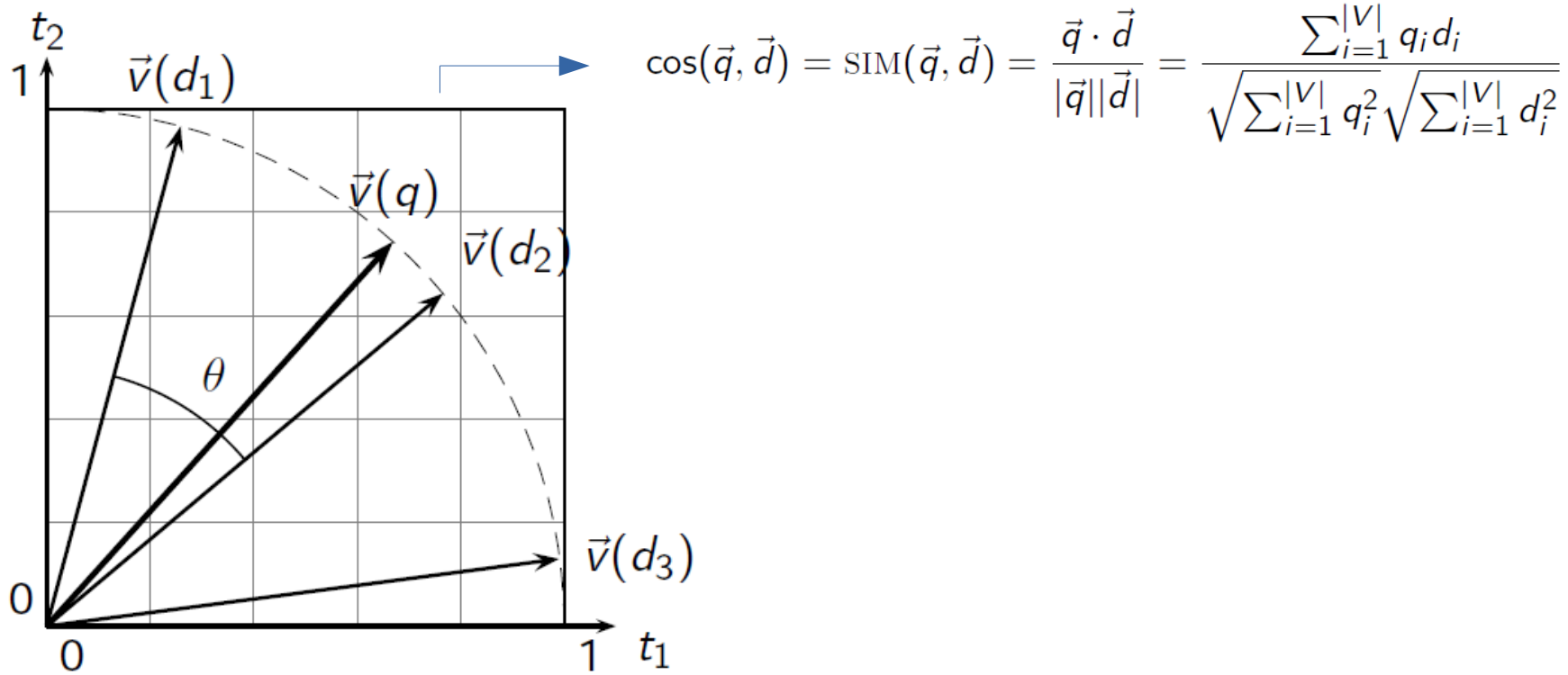
$$b \in [0, 1] , \quad k_1 > 0$$

Empírico:  $b \approx 0.75$

$$k_1 \approx 1.2$$

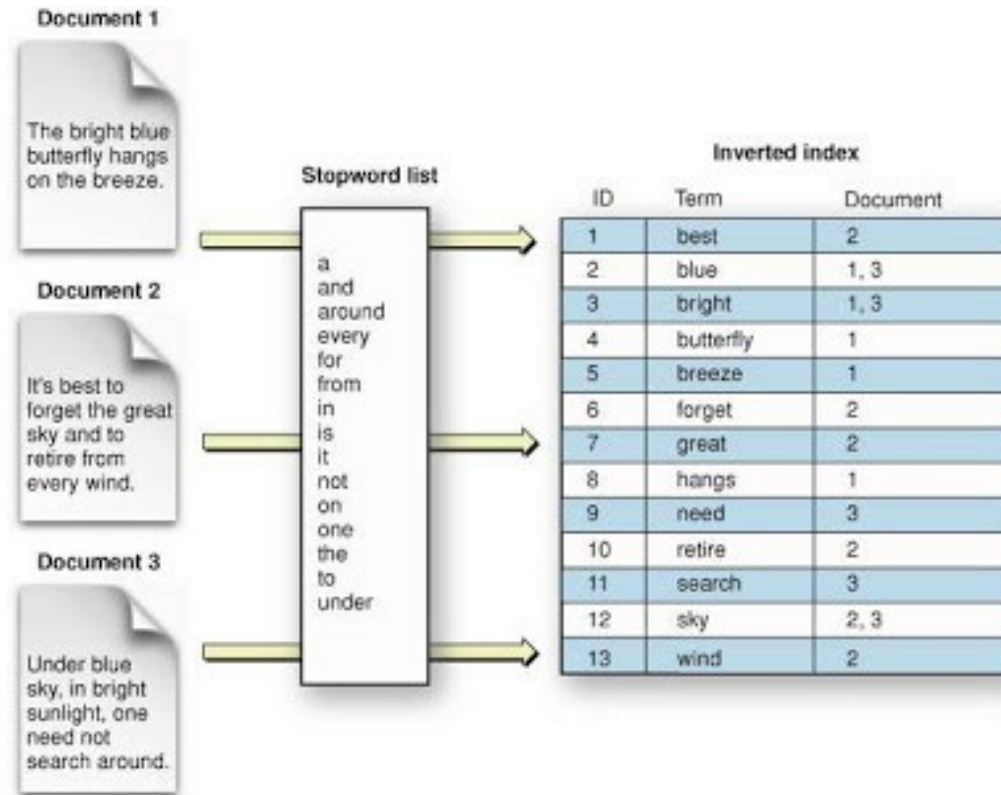
## Document ranking

Funciones de proximidad entre vectores:



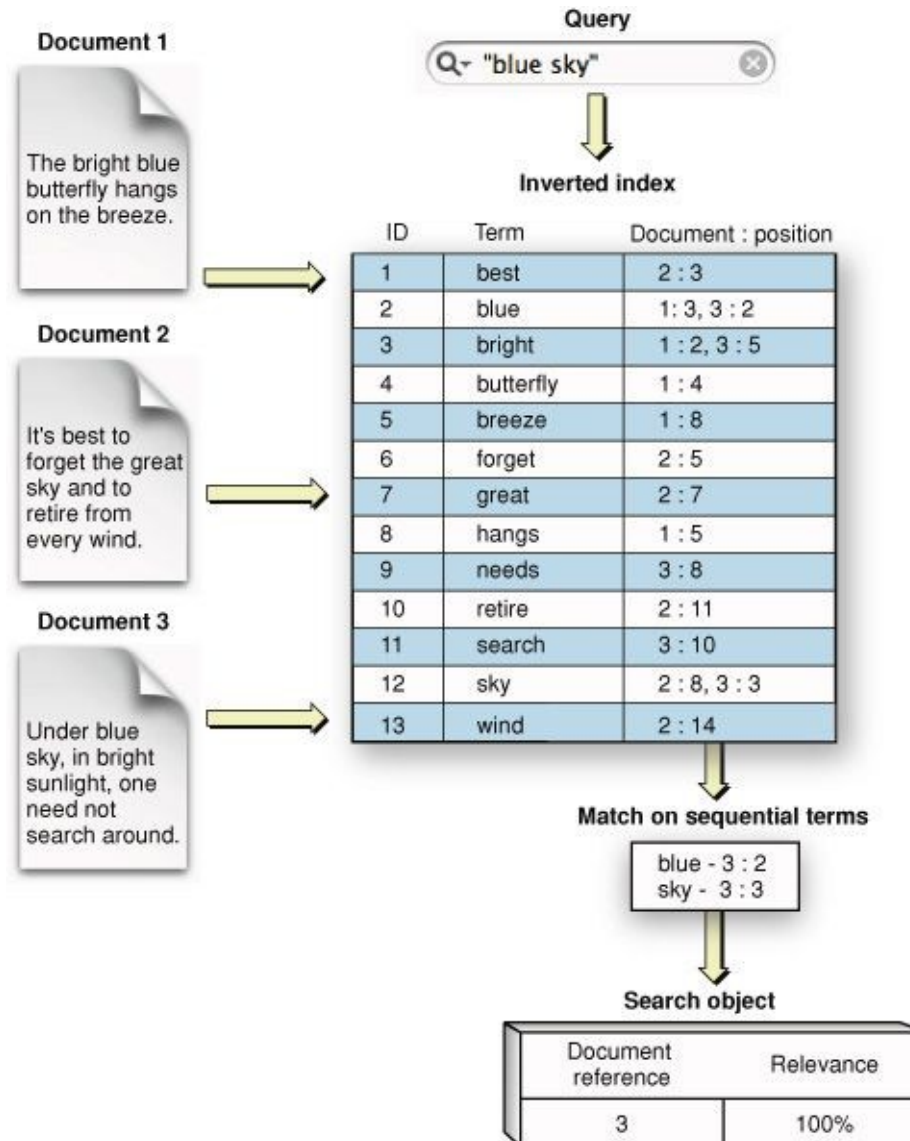
BM25 está basado en esta idea.

## Índice invertido:

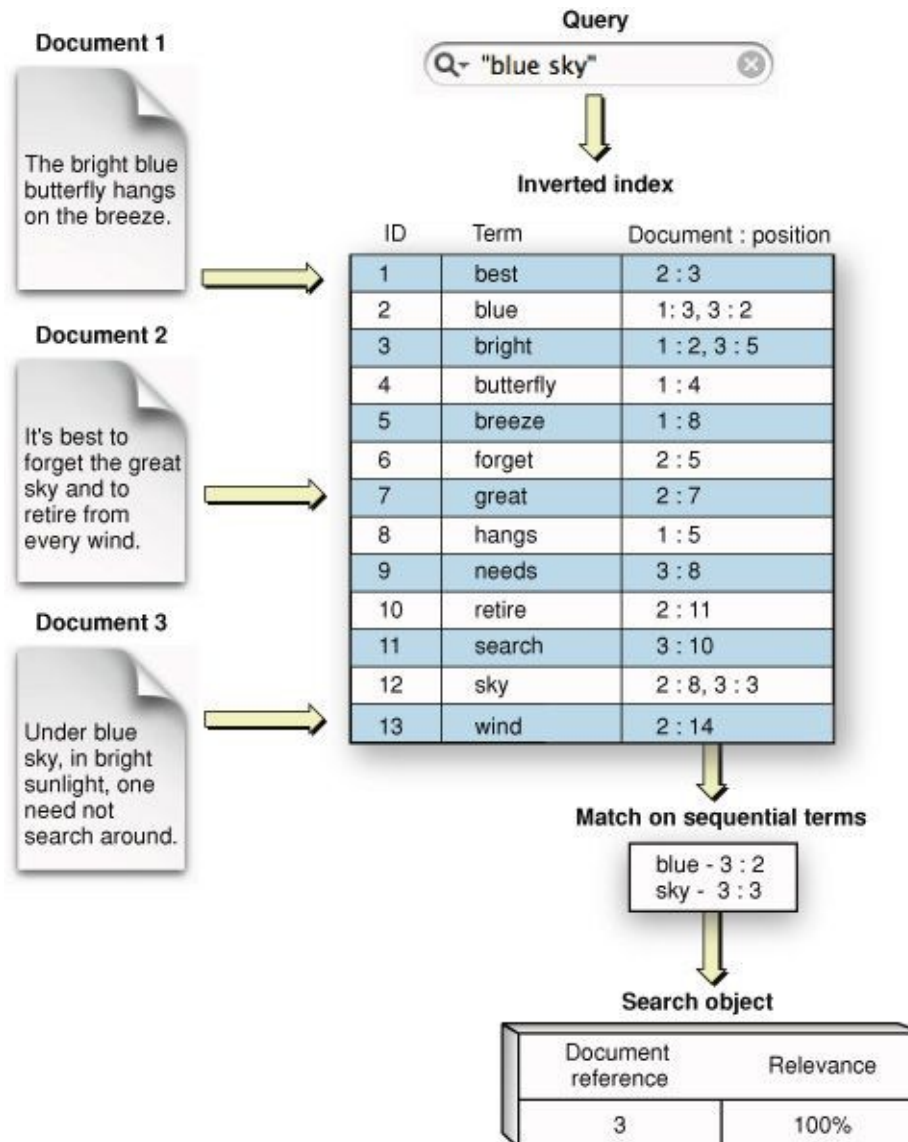




## Ranking:



## Ranking:



<https://github.com/marcelomendoza/IIC3800>