$$\text{TF-IDF} \rightarrow \text{TF} \times \text{IDF}$$

$$F() \times F() \rightarrow \text{Cantidad} \rightarrow \text{Matriz}$$
$$(\text{documento - término})$$
$$(\text{document - term})$$
$$\rightarrow \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \rightarrow \begin{matrix} \text{SVM} \\ \text{NB} \\ \text{MLP} \end{matrix}$$

Term Frequency : Cantidad de ocurrencias de una palabra $w_i$ en el documento $d_i$

document Frequency : Cantidad de documentos que contienen a la palabra $v_i$

CountVectorizer() $\rightarrow$ tFidF()   (bag-of-words)

| | and | document | first | is | one | second | the | third | this |
|---|---|---|---|---|---|---|---|---|---|
| D1 This is the first document. | **0** | $1 = \log(\frac{4}{3})$ | 1 | 1 | 0 | 0 | 1 | 0 | 1 |
| D2 This document is the second document. | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 |
| D3 And this is the third one. | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 |
| D4 Is this the first document? | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 |
| $d_{F_i}$ | 1 | 3 | 2 | 4 | 1 | 1 | 4 | 1 | 4 |
| $id_{F_i}$ | $\log(\frac{4}{1})$ | $\log(\frac{4}{3})$ | $\log(\frac{4}{2})$ | $\log(\frac{4}{4})$ | $\cdots$ | | | | |
| $t_F \times id_F$ | | | | | | | | | |

$$w_i = \underline{and}$$
$$t_F = \text{count}(t_i, d_i)$$
$$id_F = \log\left(\frac{N}{d_{F_i}}\right)$$

N: Cantidad de documentos en el corpus   [N=4]

para "and" $\rightarrow$ $t_F(\text{"and"}, D1) = 0$

$id_F\left(\frac{4}{d_{F_{and}}}\right) = \log(\frac{4}{1})$

$0 \times \log(\frac{4}{1}) = 0$

Matrices con muchos ceros $\rightarrow$ sparse matrix
$\qquad\qquad\qquad \underset{\text{sin}}{\vee} \rightarrow$ dense matrix

TFIDF   vs   TF (Bag-of-words)

ponderaciones          enfoque por
$\qquad \downarrow$          contador

Palabras:
  menos repetidas : $\uparrow$

mas repetidas: ↓

$$\text{idf}(t) = \log \frac{1+n}{1+\text{df}(t)} + 1, \quad \text{(sklearn)}$$

$$\log\left(\frac{1+4}{1+3}\right) + 1 = \log\left(\frac{5}{4}\right) + 1$$