

# Matrices terme-document

March 12, 2024

## Pondération tf-idf

Ce poids dépend de la fréquence du terme  $t$  dans le document  $d$  (tf, pour *term frequency* et du nombre des documents  $D$  qui contiennent le mot (idf, pour 'inverse doc frequency')) :

$$tf - idf_{t,d} = tf_{t,d} \cdot idf_{t,D} \quad (1)$$

où **tf** est la fréquence de  $t$  dans  $d$ , et

$$idf_{t,D} = \log \frac{|D|}{|\{d \in D : t \in d\}|} \quad (2)$$

## Sélectionner des cellules d'une matrice selon des critères

```
>>>
>>>
>>>
>>> x
array([[3, 4, 5],
       [6, 7, 8]])
>>> x[x==4]
array([4])
>>>
>>> x[x>4]
array([5, 6, 7, 8])
>>>
>>> _
```

- `x[x> 4]` retourne **les éléments de x** qui sont  $> 4$
- `np.argwhere(x>4)` retourne **les indices** des éléments de x qui sont  $> 4$
- `np.amax(x)` retourne la valeur maximale dans x
- `np.argwhere(x==np.amax(x))` retourne l'indice de l'élément avec la valeur maximale

### Exercices.

1. Nous allons construire une matrice terme-document à partir d'un extrait du corpus Reuters dans NLTK (400 documents numérotés de 0 à 399). Téléchargez les fichiers 'reuters', term-doc.py, terms.lst.
  - Les lignes de la matrice correspondent aux documents, les colonnes aux termes, les cellules au nombre d'occurrences du terme dans le document. Seul les 1000 termes figurant dans la liste terms.lst sont pris en compte. Complétez term-doc.py selon les instructions dans le fichier pour créer la matrice.
2. Complétez term-doc.py pour retourner :
  - le nombre maximum d'occurrences d'un mot dans un document observé dans le corpus
  - le nombre maximum d'occurrences d'un mot précis dans un document du corpus
  - pour un mot précis, le document dans lequel il a le plus d'occurrences
  - pour un mot précis, le nombre de documents dans lesquels il apparaît
3. Créez une deuxième matrice dans laquelle le nombre d'occurrences est remplacé par 1 si le mot apparaît dans le document, et par 0 s'il n'apparaît pas.
4. Créez une troisième matrice dans laquelle le nombre d'occurrences est remplacé par le poids tf-idf du terme dans le document.