



Recherche d'information

Prétraitements et annotation des données textuelles

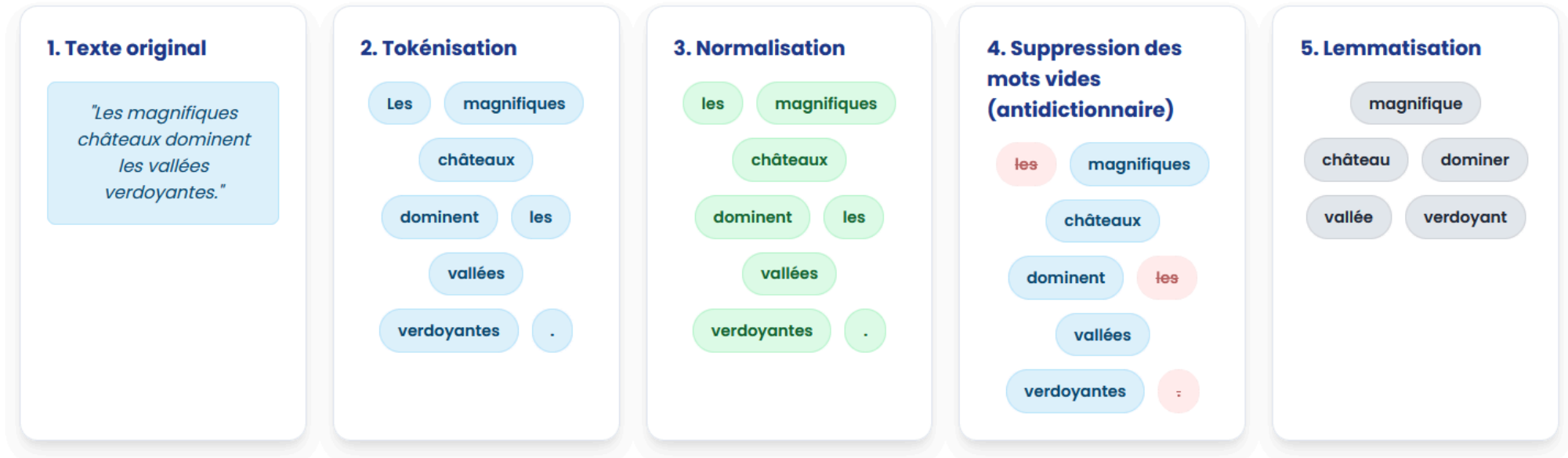
Enzo Doyen

Prétraitement des données textuelles

Avant d'indexer des documents pour de la recherche d'information, et pour améliorer l'efficacité de la recherche, il est souvent nécessaire de les traiter au préalable (en d'autres termes, opérer des « prétraitements »).

L'objectif est de passer d'un **texte brut** à un **texte lexicalement et sémantiquement riche**.

Prétraitement des données textuelles



Étapes de prétraitement des données textuelles : tokénisation, normalisation, suppression des mots vides (*stopwords*) et lemmatisation




I. Méthodes de prétraitement

Tokénisation

Original	Tokénisation simple	Tokenisation avancée
l'école	["l", "'", "école"]	["l'", "école"]
n'a	["n", "'", "a"]	["n'", "a"]
aujourd'hui	["aujourd", "'", "hui"]	["aujourd'hui"]
arrière-grand-père	["arrière", "-", "grand", "-", "père"]	["arrière-grand- père"]
Bourg-en-Bresse	["Bourg", "-", "en", "-", "Bresse"]	["Bourg-en-Bresse"]
jean.d@email.fr	["jean", ".", "d", "@", "email", ".", "fr"]	["jean.d@email.fr"]
12€50	["12", "€", "50"]	["12€50"]

Tokénisation avec *NLTK*

```
1  import nltk
2  from nltk.tokenize import word_tokenize
3  nltk.download("punkt")
4  text = "Les chiens ont l'habitude d'aboyer tous les matins."
5  tokens = word_tokenize(text, language="french")
```

 Python

Sortie: ['Les', 'chiens', 'ont', "l'habitude", 'd', 'aboyer', 'tous', 'les', 'matins', '.']

Tokénisation avec *spaCy*

```
1 import spacy
2
3 nlp = spacy.load("fr_core_news_sm")
4 text = "Les chiens ont l'habitude d'aboyer tous les matins."
5 doc = nlp(text)
6 tokens = [token.text for token in doc]
```



***spaCy* et modèles pour le français**

Dans *spaCy*, modèles de différentes tailles disponibles pour le français, avec différents corpus d'entraînement et différentes architectures.

Un modèle plus grand donne généralement de meilleures performances, mais est aussi plus lourd et plus lent à charger/utiliser.

<https://spacy.io/models/fr>

Modèles disponibles : fr_core_news_sm (petit), fr_core_news_md (moyen), fr_core_news_lg (grand), fr_core_news_trf (Transformer)

Normalisation

La **normalisation** consiste à standardiser le texte en éliminant les variations, pour ainsi faciliter le traitement ultérieur.

Normalisation

La **normalisation** consiste à standardiser le texte en éliminant les variations, pour ainsi faciliter le traitement ultérieur.

```
~ ➔ cat zoe.js
console.log('Is Zoë, Zoë?')
console.log('Zoë' === 'Zoë')
~ ➔ node zoe.js
Is Zoë, Zoë?
false
~ ➔
```

Source : <https://withblue.ink/2019/03/11/why-you-need-to-normalize-unicode-strings.html> 

Normalisation

La **normalisation** consiste à standardiser le texte en éliminant les variations, pour ainsi faciliter le traitement ultérieur.

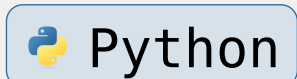
Par exemple, une normalisation voulue en français pour le mot « œuf » :

œuf = oeuf = Œuf = OEUF = ŒUF

Normalisation : casse

La normalisation au niveau de la casse se fait généralement en convertissant tout le texte en minuscules :

```
1 _ = "OEUF".lower()
```



Normalisation : espaces

```
1 _ = " oeuf ".strip()
```



```
2 _ = " oeuf".lstrip()
```

```
3 _ = "oeuf ".rstrip()
```

```
4 _ = " ".join(" un oeuf sur la table".split())
```

Normalisation : diacritiques et ligatures

Source		NFD		NFC		NFKD		NFKC
fi FB01	:	fi FB01		fi FB01		f i 0066 0069		f i 0066 0069
2 ⁵ 0032 2075	:	2 5 0032 2075		2 5 0032 2075		2 5 0032 0035		2 5 0032 0035
ḟ 1E9B 0323	:	f ̊ ̊ 017F 0323 0307		ḟ ̊ 1E9B 0323		s ̊ ̊ 0073 0323 0307		ṡ 1E69

Source : <https://www.unicode.org/reports/tr15/>

Normalisation : diacritiques et ligatures

```
1 import unicodedata
2 unicodedata.normalize('NFKC', "IJ") # 'IJ'
3 unicodedata.normalize('NFKD', "½") # '1/2'
```

 Python

```
1 from unicode import unicode
2 unicode("épée", "utf-8") # 'epee'
```

 Python

Attention

Certaines ligatures, considérées comme des lettres à part entières, ne sont pas décomposables directement (p. ex. « œ », « æ »).

Suppression de la ponctuation

```
1 import string
```

 Python

```
2
```

```
3 text = text.translate(str.maketrans("", "",  
string.punctuation))
```


Suppression des mots vides (*stopwords*)

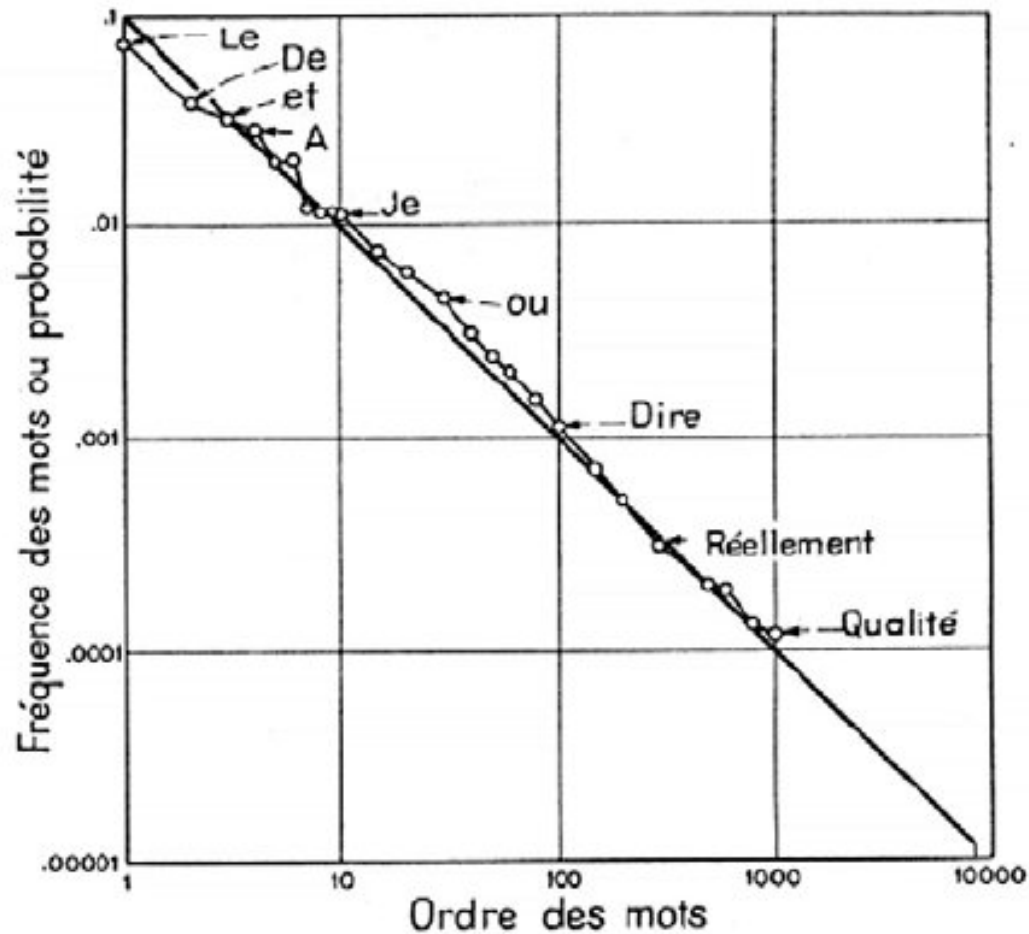
Les mots vides sont des mots qui n'apportent pas de valeur sémantique à un document, et qui peuvent être supprimés pour alléger l'indexation.

Ils incluent bien souvent les mots grammaticaux, par exemple en français :
« le », « la », « et », « à », « de », etc.

Suppression des mots vides (*stopwords*)

Lien avec la **loi de Zipf** (Zipf, 1932) : observation empirique selon laquelle, dans un corpus de texte, la fréquence d'un mot est inversement proportionnelle à son rang dans la liste des mots par fréquence.

Un mot de rang 1 (le plus fréquent) apparaîtra deux fois plus souvent qu'un mot de rang 2, trois fois plus souvent qu'un mot de rang 3, etc.



Représentation graphique de la loi de Zipf. Source : **Mandelbrot (1965)**

Suppression des mots vides (*stopwords*) avec *NLTK*

```
1 import nltk
2 from nltk.corpus import stopwords
3
4 nltk.download("stopwords")
5 english_sw = stopwords.words('english')
6 french_sw = stopwords.words('french')
```



Suppression des mots vides (*stopwords*) avec *spaCy*

```
1 import spacy
2 from spacy.lang.en import stop_words
3
4 nlp = spacy.load('en_core_web_sm')
5
6 stop_words = stop_words.STOP_WORDS
```



Suppression des mots vides (*stopwords*)

Liste en français pour spaCy : https://github.com/explosion/spaCy/blob/master/spacy/lang/fr/stop_words.py 


```
STOP_WORDS = set(
    """
a à â abord afin ah ai aie ainsi ait allaient allons
alors anterieur anterieure anterieures antérieur antérieure antérieures
apres après as assez attendu au
aupres auquel aura auraient aurait auront
aussi autre autrement autres autrui aux auxquelles auxquels avaient
avais avait avant avec avoir avons ayant

bas basee bat
```

Google


All Images Videos Short videos Web News Books More ▾

Words List Dictionary

 GitHub
[https://github.com › stopwords-iso › stopwords-fr](https://github.com/stopwords-iso/stopwords-fr) ⋮


stopwords-iso/stopwords-fr: French stopwords collection

The most comprehensive collection of **stopwords** for the **french** language. A multiple language collection is also available.

 GitHub
[https://github.com › gillesbastin › french_stopwords](https://github.com/gillesbastin/french_stopwords) ⋮

gillesbastin/french_stopwords: A list of stopwords to be ...

This repository contains a carefully curated list of **stopwords** for preprocessing **French** texts, specifically designed for text mining tasks (**FRENCH_STOPWORDS**).

 GitHub
[https://github.com › stopwords-json](https://github.com/stopwords-json) ⋮

Stopwords for 50 languages in JSON format

Lemmatisation

La lemmatisation est le processus de réduction des mots à leur forme canonique (leur « lemme »), ce qui permet de regrouper les différentes formes d'un mot.

Lemmatisation avec *spaCy*

```
1 import spacy
2
3 nlp = spacy.load("fr_core_news_sm")
4 text = "Les chiens ont l'habitude d'aboyer tous les matins."
5 lemmas = [token.lemma_ for token in nlp(text)]
```



Sortie: ['le', 'chien', 'avoir', 'le', 'habitude', 'de', 'aboyer', 'tout', 'le', 'matin', '.']

Lemmatisation avec *treetaggerwrapper*

```
1 import treetaggerwrapper
```

 Python

```
2
```

```
3 tagger = treetaggerwrapper.TreeTagger(TAGLANG='fr')
```

```
4 tags =  
treetaggerwrapper.make_tags(tagger.tag_text(text))
```

```
5
```

```
6 lemmas = [t.lemma for t in tags if t.lemma is not  
None]
```

Autres types de prétraitements

- ♦ corrections de fautes d'orthographe ;
- ♦ correction phonétique ;
- ♦ *chunking* (découpage en phrases ou en segments, notamment pour les systèmes de *RAG*) ;
- ♦ suppression d'informations sensibles (par exemple, pour des raisons de confidentialité) ;
- ♦ ...



II. Annotation automatique

Annotation automatique

Il est possible d'utiliser des techniques du domaine de l'**extraction d'informations** pour annoter automatiquement des documents avec des informations complémentaires.

Ces techniques peuvent permettre de récupérer des mots-clés du document (extraction de mots-clés, *keyword extraction*), des entités nommées (reconnaissance d'entités nommées, *named entity recognition*, NER), des relations entre entités (extraction des relations, *relationship extraction*), etc.

Les informations additionnelles obtenues grâce à ces techniques peuvent être utilisées pour améliorer la recherche et la visibilité des résultats, ou bien pour filtrer/regrouper des documents, par exemple.

Annotation automatique : extraction de mots-clés

Exploitation de différentes caractéristiques linguistiques pour trouver les mots les plus pertinents dans le texte : casse, position, fréquence au sein du document, fréquence au niveau des phrases, etc.

Exemples d'implémentation en Python :

- ♦ Yake (**Campos et al. (2018)** ; notebook disponible sur Moodle)
- ♦ pke (**Boudin (2016)** ; notebook disponible sur Moodle)

Annotation automatique : extraction de mots-clés

```
1 import yake
2 kw_extractor = yake.KeywordExtractor()
3 keywords = kw_extractor.extract_keywords(text)
```

 Python

```
1 ('google', 0.026580863364597897) # score
   plus bas = mot plus pertinent
2 ('kaggle', 0.0289005976239829)
3 ('ceo anthony goldbloom', 0.029946071606210194)
4 ('san francisco', 0.048810837074825336)
```

 Python

Annotation automatique : extraction de mots-clés

Utilisation de plongements lexicaux pour de meilleures performances, par exemple KeyBERT (basé sur sentence-transformers) avec un modèle multilingue [?].

Annotation automatique : extraction de mots-clés

```
1 from keybert import KeyBERT
```

 Python

```
2 kw_model = KeyBERT()
```

```
3 keywords = kw_model.extract_keywords(doc)
```

```
1 [ ('learning', 0.4604), # score plus haut =  
   mot plus pertinent
```

 Python

```
2 ( 'model', 0.4588),
```

```
3 ( 'algorithm', 0.4556),
```

```
4 ( 'training', 0.4487)]
```

Annotation automatique : extraction de mots-clés

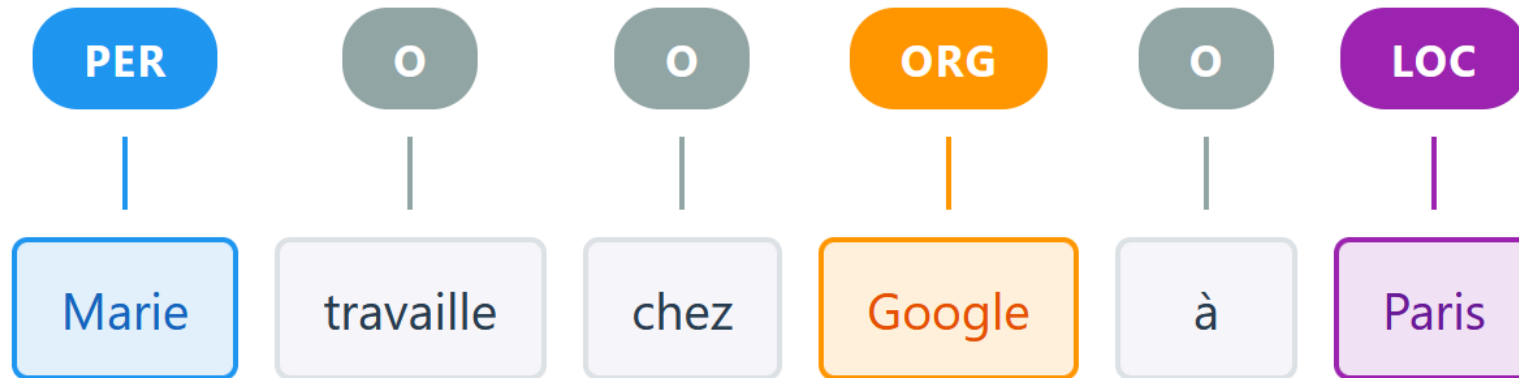
Quelques exemples d'options personnalisées avec *Yake* (voir la [documentation officielle](#)) :

```
1 custom_kw_extractor =  
  yake.KeywordExtractor(  
2     lan="en",           # langue  
3     n=3,                # ngram  
4     windowsSize=1,      # fenêtre de contexte  
5     top=10,             # nb mots-clés à extraire  
6 )
```



Annotation automatique : reconnaissance d'entités nommées

La reconnaissance d'entités nommées (NER) consiste à identifier et classifier des entités nommées (noms de personnes, d'organisations, de lieux, dates, etc.) dans un texte.



Annotation automatique : reconnaissance d'entités nommées

La reconnaissance d'entités nommées (NER) consiste à identifier et classifier des entités nommées (noms de personnes, d'organisations, de lieux, dates, etc.) dans un texte.

Diverses options disponibles :

- ♦ spaCy (notebook disponible sur Moodle) ;
- ♦ Flair (modèle disponible sur [HuggingFace](#)) ;
- ♦ DistilCamemBERT-NER (modèle disponible sur [HuggingFace](#) et basé sur CamemBERT).

Annotation automatique : reconnaissance d'entités nommées

Utilisation :

```
1 from transformers import pipeline
2 ner = pipeline(
3     task='ner',
4     model="cmarkea/distilcamembert-base-ner",
5     tokenizer="cmarkea/distilcamembert-base-ner",
6     aggregation_strategy="simple" # entités multi-
    tokens regroupées
7 )
```



Annotation automatique : reconnaissance d'entités nommées

Exemple de sortie :

```
1 [{ 'entity_group': 'ORG',  
2   'score': 0.9974479,  
3   'word': 'Crédit Mutuel Arkéa',  
4   'start': 3,  
5   'end': 22},
```



Bibliographie

- Boudin, F. (décembre 2016). Pke: An Open Source Python-Based Keyphrase Extraction Toolkit. In H. Watanabe (éd.), *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: System Demonstrations: Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: System Demonstrations*.
- Campos, R., Mangaravite, V., Pasquali, A., Jorge, A. M., Nunes, C., et Jatowt, A. (2018). A Text Feature Based Automatic Keyword Extraction Method for Single Documents. In G. Pasi, B. Piwowarski, L. Azzopardi, et A. Hanbury (Éds.), *Advances in Information Retrieval: Advances in Information Retrieval*. [10.1007/978-3-319-76941-7_63](#)
- Mandelbrot, B. B. (1965). Information Theory and Psycholinguistics. In B. B. Wolman et E. Nagel (éds.), *Scientific Psychology: Scientific Psychology*. Basic Books.
- Manning, C. D., Raghavan, P., et Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press.
- Mitra, B. (2018). *An Introduction to Neural Information Retrieval* (Numéro v.41). Now Publishers.
- Zhai, C., et Massung, S. (juin 2016). *Text Data Management and Analysis: A Practical Introduction to Information Retrieval and Text Mining*. Association for Computing Machinery and Morgan & Claypool. [10.1145/2915031](#)
- Zipf, G. K. (1932). *Selected Studies of the Principle of Relative Frequency in Language* (p. 57). Harvard Univ. Press. [10.4159/harvard.9780674434929](#)