

# Démonstration spaCy

## Projet tutoré FIE5 ISIS 2023

AYMARD Marion | JOLIVET Anaëlle | BOURDARIAS Cléa

Entrée [1]:

```
1 # Importe spaCy
2 import spacy
3
4 # Crée un objet nlp français vide
5 nlp = spacy.blank("fr")
```

Le paquet "fr\_core\_news\_sm" est un petit pipeline pour la langue française qui propose toutes les fonctionnalités de base et qui a été entraîné sur des textes d'actualité. Ici on utilise le moyen (**fr\_core\_news\_md**) car le petit ne supporte pas certaines fonctions que l'on va utiliser.

Entrée [2]:

```
1 # Charge le petit pipeline français
2 nlp = spacy.load("fr_core_news_md")
```

On utilise la phrase **"Je me suis fait mal au poignet hier après-midi."**.

Entrée [3]:

```
1 # Crée en traitant une chaîne de caractères avec l'objet nlp
2 doc = nlp("Je me suis fait mal au poignet hier après-midi.")
```

## Prédire les types de mot dans la phrase

Entrée [4]:

```
1 # Itère sur les tokens
2 for token in doc:
3     # Affiche le texte et l'étiquette de partie de discours prédite
4     print(token.text, token.pos_)
```

```
Je PRON
me PRON
suis AUX
fait VERB
mal ADV
au ADP
poignet NOUN
hier ADV
après-midi NOUN
. PUNCT
```

## Prédire comment les mots sont reliés

Entrée [5]:

```
1 for token in doc:
2     print(token.text, token.pos_, token.head.text)
```

Je PRON fait  
me PRON fait  
suis AUX fait  
fait VERB fait  
mal ADV fait  
au ADP poignet  
poignet NOUN fait  
hier ADV fait  
après-midi NOUN fait  
. PUNCT fait

## Prédire des similitudes entre les phrases

On prend les phrases "J'ai mal aux dents." et "J'ai une douleur dentaire.". On peut déterminer un taux de similitude entre ces phrases.

Entrée [6]:

```
1 # Compare deux documents
2 doc1 = nlp("J'ai mal aux dents.")
3 doc2 = nlp("J'ai une douleur dentaire.")
4 print(doc1.similarity(doc2))
```

0.8645883336750065

On prend les phrases "J'ai mal aux dents." et "J'ai une douleur aux pieds."

Entrée [7]:

```
1 # Compare deux documents
2 doc1 = nlp("J'ai mal aux dents.")
3 doc2 = nlp("J'ai une douleur aux pieds.")
4 print(doc1.similarity(doc2))
```

0.927480481745581

On prend les phrases "Je n'ai pas mal aux dents." et "J'ai mal aux dents."

Entrée [8]:

```
1 # Compare deux documents
2 doc1 = nlp("Je n'ai pas mal aux dents.")
3 doc2 = nlp("J'ai mal aux dents.")
4 print(doc1.similarity(doc2))
```

0.7649048822216306

Cette fonction n'est **pas** forcément **pertinante** car on remarque que les phrases "J'ai mal aux dents." et "J'ai une douleur aux pieds." sont plus similaires que les phrases "J'ai mal aux dents." et "J'ai une douleur dentaire.". De plus, on remarque un taux de similitude important entre les phrases "Je n'ai pas mal aux dents." et "J'ai mal aux dents." alors qu'elles sont à caractère opposé.

## Recupérer les mots qui nous intéressent

Toujours avec notre phrase "**Je me suis fait mal au poignet hier après-midi.**", on peut déterminer les mots qui nous semblent importants.

```
Entrée [9]: 1 # Initialise avec Le vocabulaire partagé
2 from spacy.matcher import Matcher
3 matcher = Matcher(nlp.vocab)
```

```
Entrée [10]: 1 doc = nlp("Je me suis fait mal au poignet hier après-midi.")
```

```
Entrée [11]: 1 # Écris un motif qui recherche une forme de "mal" suivie d'un nom
2 pattern = [{"LEMMA": "mal"}, {"POS": "ADP"}, {"POS": "DET", "OP": "?"}]
3
4 # Ajoute Le motif au matcher et applique Le matcher au doc
5 matcher.add("DOULEUR", [pattern])
6 matches = matcher(doc)
7 print("Nombre de correspondances trouvées :", len(matches))
8
9 # Itère sur Les correspondances et affiche La portion de texte
10 for match_id, start, end in matches:
11     print("Correspondance trouvée :", doc[start:end].text)
```

```
Nombre de correspondances trouvées : 1
Correspondance trouvée : mal au poignet
```

Le format de ce pattern permet de trouver le même résultat même si la phrase est légèrement différente. Par exemple, essayons avec la phrase "**Je me suis fait mal à la main hier après-midi.**".

```
Entrée [12]: 1 doc = nlp("Je me suis fait mal à la main hier après-midi.")
2
3 matches = matcher(doc)
4 print("Nombre de correspondances trouvées :", len(matches))
5
6 for match_id, start, end in matches:
7     print("Correspondance trouvée :", doc[start:end].text)
```

```
Nombre de correspondances trouvées : 1
Correspondance trouvée : mal à la main
```

On remarque qu'il repère le bout de phrase "**mal à la main**" même si "**mal**" et "**main**" sont séparés de deux mots et non d'un seul comme dans "**mal au poignet**".

# Récupérer les mots importants regroupés dans un fichier JSON

Il est possible de comparer les mots de notre doc avec ceux d'une liste d'un fichier json. Ici on récupère les mots que nous voulons trouver dans le texte depuis le fichier json **"vocabulaire.json"**.

Entrée [22]:

```
1 import json
2
3 with open("vocabulaire.json", encoding="utf8") as f:
4     VOCABULAIRE = json.loads(f.read())
5
6 print(VOCABULAIRE)
```

```
['mal', 'poignet', 'hier', 'aïe', 'jambe', 'voiture']
```

Ensuite, dans la phrases **"Je me suis fait mal au poignet hier après-midi."** on ne recupère que les mots intéressants, c'est à dire, ceux qui sont dans le fichier json.

Entrée [23]:

```
1 nlp = spacy.blank("fr")
2 doc = nlp("Je me suis fait mal au poignet hier après-midi.")
3
4 # Importe Le PhraseMatcher et initialise-Le
5 from spacy.matcher import PhraseMatcher
6
7 matcher = PhraseMatcher(nlp.vocab)
8
9 # Crée des motifs objets Doc et ajoute-Les au matcher
10 # C'est La version rapide de : [nlp(country) for country in COUNTRIES]
11 patterns = list(nlp.pipe(VOCABULAIRE))
12 matcher.add("VOCABULAIRE", patterns)
13
14 # Appelle Le matcher sur Le document de test et affiche Le résultat
15 matches = matcher(doc)
16 print([doc[start:end] for match_id, start, end in matches])
```

```
[mal, poignet, hier]
```

On remarque que les mots sélectionnés ne sont que ceux qui sont listés dans le fichier **"vocabulaire.json"**.