

MASTER TRAITEMENT AUTOMATIQUE DES LANGUES

Ingénierie des connaissances

Rapport

Mise à disposition d'un corpus annoté pour étudier la
variation sociolinguistique diastratique

Étudiants :

Kenza Ahmia
Florian Jacquot
Tiffany Nguyen
Shami Thirion Sen

Année universitaire 2023/2024

Table des matières

1	Introduction	2
2	Présentation des corpus	3
2.1	Corpus ESLO	3
2.2	Corpus 89milSMS	3
3	Réflexion et parcours suivi.	4
3.1	Temps verbaux	4
3.2	Négation	5
3.3	Le /y/ et schwa	5
3.4	Les formats de données d'entrée et de sortie	6
4	Méthodes utilisées	6
5	Résultat et retour d'expérience	10
5.1	Observation des résultats	10
5.2	Sorties	11
6	Conclusion	11
A	Scripts Python utilisés	13
A.1	Script pour l'annotation du corpus SMS et faire la sortie XML	13
A.2	Script pour l'annotation du corpus ESLO	18
A.3	Script pour l'extraction du corpus sms	22
A.4	Script pour l'extraction du corpus ESLO	25
A.5	Script pour la création d'un fichier niveau d'études pour ESLO	26
A.6	Script pour transformer le fichier ESLO annoté en format XML	26

1 Introduction

Contribution : Florian

Le présent rapport a pour objet de présenter le travail effectué dans le cadre d'un projet de recherche visant à étudier la variation sociolinguistique diastratique dans deux types de données : des données issues de l'oral et de la communication médiée par les réseaux. Notre équipe est chargée de contribuer à l'avancement des connaissances dans ce domaine. Tous les scripts, corpus brut et annotés se trouvent sur le git du projet ¹.

Les variables étudiées dans le cadre de ce projet sont les temps verbaux, la négation complète et incomplète, ainsi que l'absence ou la présence de /y/ dans « tu » et du schwa dans les clitiques. Pour mener à bien cette étude, nous avons suivi un processus en six étapes :

1. faire l'inventaire des sources de connaissances disponibles (corpus)
2. acquisition des connaissances (corpus)
3. modélisation formelle des connaissances acquises
4. modélisation formelle des connaissances à ajouter (modèles d'annotation)
5. ajout de nouvelles connaissances (annotation automatique)
6. formalisation des connaissances obtenues (création d'une nouvelle version du corpus)

Dans ce rapport, nous détaillerons chacune de ces étapes et présenterons les résultats obtenus. Nous expliquerons également les choix méthodologiques que nous avons faits et les éventuelles difficultés rencontrées. Enfin, nous présenterons les nouvelles versions des corpus ESLO et 88milSMS, que nous avons généré dans le cadre de ce projet.

Répartition des tâches

Note : les contributions indiquées sous chaque titre de partie du rapport ne concernent que les contributions à la rédaction et n'ont rien à voir avec le travail fait en amont. Pour ce travail en particulier, la répartition fut plutôt la suivante :

- Extraction des connaissances du corpus ESLO : **Tifanny**
- Extraction des connaissances du corpus 88milSMS : **Shami**
- Annotation des corpus : **Kenza et Florian**
- Sortie XML : **Kenza et Tifanny**
- Exportation TXM : **Kenza et Florian**

Mais pour chacune de ces étapes, tous les membres du groupe ont participé à la réflexion.

1. https://github.com/sham-glam/Ingenierie_Connaissance.git

2 Présentation des corpus

Contribution : Florian

Le travail a commencé par une recherche de deux corpus de natures différentes pour avoir des données contrastées. Les corpus retenus furent les corpus ESLO et 88milSMS. Les corpus étant plutôt compliqué à récupérer, notre professeur Loïc Liegeois nous les a fournis.

2.1 Corpus ESLO

Le corpus ESLO – Enquête SocioLinguistique à Orléans – est un corpus de données linguistiques orales retranscrites, recueillies dans le cadre d’une enquête sociolinguistique menée à Orléans, en France. Il a été créé par une équipe de chercheurs en linguistique et en sociolinguistique et est destiné à être utilisé pour l’étude de la variation linguistique dans la ville d’Orléans. Le corpus a été constitué en deux périodes. Un premier travail a été fait entre 1969 et 1974, corpus ESLO1 et un second en 2014, ESLO2. Étant donné que le projet se concentre sur la période contemporaine, seul le corpus ESLO2 est utilisé ici.

Le corpus comprend quatre-vingt quatre fichiers XML, correspondants à quatre-vingt quatre enregistrements audio, ainsi qu’un fichier `tsv` avec des métadonnées supplémentaires. Dans les fichiers XML, chaque unité de parole possède un attribut indiquant l’identifiant du locuteur. Par exemple, l’unité donnée ci-dessous nous indique que le locuteur avec l’identifiant BV1 a énoncé « d’accord » entre 600.637s et 600.862s de l’enregistrement.

```
<u who="BV1" start="600.637" end="600.862">d'accord</u>
```

De son côté, le fichier `tsv` contient des métadonnées sur les locuteurs tels que leur niveau d’étude, ou leur situation familiale. Pour les besoins du projet, seuls les énoncés, les identifiants et les niveaux d’étude furent conservés (voir section ??).

2.2 Corpus 89milSMS

Le corpus 88milSMS est un corpus de données linguistiques écrites, constitué de plus de 88 000 SMS en français, collectés en 2011 pour dans le cadre du projet sud4science LR². Lors de la collecte, des données sociolinguistiques ont été recueillies sur les auteurs des SMS. Le corpus se présente sous la forme d’un unique fichier XML au format de balisage TEI, contenant toutes les informations recueillies. Le format d’encodage TEI permet de faire des relations entre les différents éléments contenus dans le fichier. Ainsi, ce dernier présente plusieurs sections. Une première, présente des métadonnées avec toutes les informations sur les auteurs des SMS, dont l’identifiant les représentant, le niveau d’éducation, l’âge ou encore

2. <http://88milsms.huma-num.fr/>

Corpus	Nombre tokens
ESLO	1 131 986
88milSMS	932 471

TABLE 1 – Nombre de tokens des corpus ESLO2 et 88milSMS utilisés.

le type de téléphone utilisé pour écrire les SMS. Une seconde section présente les SMS qui sont associés aux identifiants des auteurs. Comme pour le corpus ESLO2, nous avons extrait de ce document les énoncés des SMS, les identifiants et les niveaux d'études.

Contrairement au corpus ESLO2, qui est constitué de données orales retranscrites avec des règles définies, le corpus 88milSMS est constitué de données écrites, avec toutes les particularités orthographiques et syntaxiques que cela implique. Par exemple, les SMS peuvent contenir des abréviations qui ne seront pas présentes dans les retranscriptions. Ces particularités peuvent être intéressantes à étudier en elles-mêmes, mais elles peuvent aussi compliquer la tâche d'annotation. Le corpus ESLO2, quant à lui, peut également amener son lot de difficultés pour l'annotation, avec les erreurs de syntaxes qui apparaissent à l'oral du fait des répétitions, des hésitations, des coupures de phrases etc.

Nous présentons dans le tableau 1 le nombre de tokens totaux de chacun de ces deux corpus.

3 Réflexion et parcours suivi.

Contribution : Tiffany, Florian

Le projet, comme indiqué précédemment, nous demandait d'identifier dans les deux corpus les temps verbaux, la négation complète et incomplète, l'absence ou la présence de /y/ dans « tu » et du schwa.

Avant de nous mettre à faire des scripts, nous avons pris un temps de réflexion afin de pouvoir coordonner le travail sur les deux corpus et également trouver comment faire pour identifier tout ce qui nous a été demandé.

3.1 Temps verbaux

Pour les temps verbaux, il nous fallait un outil qui permettait d'identifier les verbes mais aussi de nous donner son temps. Le travail étant assez clair, il n'y avait rien de spécial à identifier si ce n'est que l'outil que nous allions utiliser. Au début, nous nous sommes penchés sur TreeTagger qui semblait pouvoir faire le travail. Cependant, ayant rencontrés de nombreux problèmes à son installation, nous avons opté pour Spacy. Dans le cadre de ce travail, nous avons choisi de travailler avec le modèle français le plus léger de Spacy mais il est tout à fait possible de changer ce paramètre dans le script d'annotation A.1 pour obtenir un meilleur résultat.

Ainsi, pour obtenir une annotation sur les temps verbaux, nous avons utilisé l’annotation en partie du discours de Spacy qui nous permet d’avoir ces informations.

3.2 Négation

Il y a deux types de négation : la complète et l’incomplète. La différence entre les deux est le « ne » ou « n’ » qui n’est pas présent dans la négation incomplète. Notre professeur Loïc Liegeois nous avait déjà identifié plusieurs tournures de négations :

- Il (ne) ... pas.
- Il (ne) ... point.
- Il (ne) ... rien.
- Il (ne) ... jamais.
- etc...

Il nous fallait juste un moyen de les détecter sans les lister un par un. Pour la négation complète, nous avons remarqué que dans la plupart des cas, lors d’une négation complète, le verbe est entouré par deux adverbes : « ne » et « pas » par exemple. Nous avons donc décidé d’utiliser l’annotation en partie du discours déjà implémentée pour la recherche des temps verbaux afin de pouvoir identifier la négation complète plus facilement.

Cette annotation, qui aurait pu se faire avec TreeTagger, a finalement été faite avec Spacy.

3.3 Le /y/ et schwa

L’absence ou la présence du /y/ dont nous parlons ici est celle qui accompagne le « tu » : « Tu as pris » qui se transforme en « T’as pris ». Quant au schwa, c’est l’absence du « e » dans les clitiques « ce, de, je, le, me, ne, que, se, te » dont la liste nous a été donnée par notre professeur.

Aucun outil spécifique ne permet de faire ce type d’annotation, nous avons donc défini notre propre méthode pour la réaliser. Il est à noter que le corpus ESLO ne présente pas d’absence de schwa du fait que les retranscriptions ont ignoré cette information linguistique par simplicité. L’annotation sur les schwa n’est donc présente que sur le corpus 88milSMS.

La méthode retenue ici pour effectuer cette annotation est l’utilisation de motifs – expressions régulières. En effet, ces deux variables possèdent dans une grande majorité des cas, des similarités de formes qui ne se croisent pas entre les deux cas, ce qui permet de distinguer les deux variables.

Dans le cas de l’absence de /y/ dans le « tu », nous avons remarqué que, dans la grande majorité des cas, la modification se produit – dans ce corpus – lorsque le mot suivant « tu » commence par une voyelle comme « tu as » ou « tu es » qui deviennent « t’as » (ou « tas ») et « t’es » (ou « tes ») respectivement. Le motif est donc « t’ + mot commençant par une voyelle ». Cependant, cette condition est trop large car la forme « t’es » existe aussi dans

la langue – par exemple, dans la phrase « Tu t’es perdue ? ». Nous avons donc ajouté une condition en mettant à profit notre utilisation de Spacy qui nous permet de récupérer la relation de dépendance des tokens par rapport à la tête. Ici, puisque nous sommes dans le cas d’une relation sujet, le motif utilisé est plutôt « t’ de dépendance sujet + mot commençant par une voyelle ».

3.4 Les formats de données d’entrée et de sortie

Maintenant que les outils d’identification des connaissances sont définis 3, il nous faut définir les formats de données que nous utilisons et expliquer les raisons derrière les choix réalisés.

Pour commencer, comme présenté dans la section sur les corpus 2, les données d’entrée étaient en format XML – et csv pour le corpus ESLO. Ce format étant universel, simple et à la base de beaucoup d’applications gérant des données, nous avons choisi de viser un format de sortie XML. De plus, devant mettre à disposition le corpus annoté sur un outil facilement utilisable par les futurs utilisateurs, ce format est idéal, notamment pour l’outil de textométrie TXM – qui sera l’outil de formalisation des connaissances obtenues de ce projet.

4 Méthodes utilisées

Contribution : Kenza

Comme indiqué précédemment, après des réflexions approfondies sur les outils à utiliser, les formats d’entrée et de sortie, et la manière de modéliser nos connaissances pour traiter le corpus, nous avons commencé par récupérer les deux corpus depuis leurs sources.

Nous avons soigneusement sélectionné les fichiers XML de la manière suivante : pour le corpus SMS, il n’y avait qu’un seul fichier contenant à la fois les métadonnées et les énoncés (ID des locuteurs et niveaux d’étude). Pour le corpus ESLO, en plus des nombreux fichiers d’énoncés qui contiennent uniquement les ID des locuteurs, nous avons récupéré un fichier de métadonnées pour compléter l’information sur le niveau d’étude des locuteurs.

Tous ces fichiers ont été stockés dans un dépôt Git du projet, accessible à tous les membres de l’équipe.

1. Préparation des données d’entrée

La préparation des données a été réalisée grâce à des scripts Python, dont le code est disponible en annexe.

- (a) Extraction des fichiers textes : Pour chaque corpus, nous avons extrait trois fichiers : un fichier texte contenant tous les énoncés avec un énoncé par ligne. Un deuxième fichier avec les identifiants des locuteurs dans le même ordre. Enfin, nous

avons récupéré les fichiers des niveaux d'études des locuteurs et leurs id correspondants.

- (b) Correspondance des données : Ces fichiers représentent la deuxième version du corpus que nous prenons en entrée de notre script d'annotation et qui nous permettent de faire correspondre les énoncées à annoter avec les locuteurs ainsi que leurs niveaux d'étude.

2. Annotation des données

Les données ont été annotés à l'aide d'un script python utilisant l'outil spaCy³, une bibliothèque gratuite et open source publiée sous la licence MIT pour le traitement naturel du langage. Les annotations incluent les parties du discours (POS) incluant les temps verbaux, les négations complètes, les absences **y** et la présence de **schwa**. Le processus s'est déroulée en plusieurs étapes détaillées ci-dessous :

(a) Chargement des données

Les fichiers textes contenant les énoncées du corpus ont été importés dans le script Python à l'aide de la bibliothèque de manipulation de fichiers **sys**. Une fonction de lecture a été utilisé pour extraire les énoncés, les identifiants des locuteurs ainsi que les niveaux d'études, et ont été stockés dans des variables pour une utilisation ultérieure.

(b) Traitement des données

Chaque énoncé a été analysé avec spaCy pour extraire les tokens et leurs parties du discours. La bibliothèque a été principalement utilisé pour annoter les temps verbaux grâce à sa fonction **token.morph**. Pour détecter les négations complètes, ainsi que les absences ou présences de 'y' et les 'schwa', nous avons opté pour l'utilisation des expressions régulières (**regex**) en complément car elles permettent de spécifier de manière précise les éléments que nous souhaitons identifier.

(c) Identification et annotation des données :

- **Détection des temps verbaux** : comme spécifié dans la partie Réflexion et parcours suivi, la détection des temps verbaux a été effectué avec la fonction **token.morph** de spaCy.
- **Détection des négations complètes** : pour celle-ci, nous avons d'abord élaboré deux listes : une comprenant toutes les formes de négation couramment utilisées, telles que "ne", "n'", et "n", et une autre pour tous les adverbes de négation. Cela nous a permis de couvrir un large éventail de contextes où la négation peut apparaître. Ensuite, en utilisant des expressions régulières, nous avons parcouru chaque énoncé pour identifier les occurrences de ces formes de négation. Les regex nous ont permis de rechercher des motifs spécifiques dans le texte, comme la présence de "ne" suivi immédiatement d'un adverbe de

3. <https://spacy.io/usage>

négation tel que "pas", "rien", ou "jamais". Une fois qu'une négation complète a été identifiée, elle a été annotée dans notre système d'annotation sous forme de booléen pour indiquer sa présence.

- **Détection de l'absence de y** : dans les contractions de "tu" en "t'", dans les données, nous utilisons une fois de plus la bibliothèque **regex** pour effectuer des correspondances d'expressions régulières sur le texte. Le code parcourt chaque token (mot ou caractère) dans le document. Si le token est "tu" et qu'il est suivi d'une voyelle, le code considère que le "y" est présent et ajoute la valeur **False** à la liste **Yabsent** dans le dictionnaire d'annotations. Si le token correspond à "t'" et est suivi d'une voyelle, le code considère que le "y" est absent et ajoute la valeur **True** à la liste.
- **Détection de l'absence du schwa** : Pour détecter l'absence du **schwa**, le script utilise la même bibliothèque que pour les détections précédentes et cela en parcourant chaque token (mot ou caractère) dans le corpus. Si le token correspond à un clitique se terminant par une apostrophe (par exemple "c'", "d'", "j'", "l'", "m'", "n'", "q'", "s'", "t'") et qu'il est suivi d'un caractère qui n'est pas une voyelle ou "h", le code considère que le **schwa** est absent et ajoute la valeur **True** à la liste d'annotations. De même, si le token est "j" et qu'il est suivi d'un caractère qui n'est pas "v", "f", "t", "p", "s", "d", "c", le code considère que le **schwa** est absent et ajoute également la valeur **True**.

3. Génération des données en sortie

Pour la génération de données en sortie nous avons opté pour un format dictionnaire de dictionnaire pour les résultats d'annotations que nous avons ensuite transformé en fichier XML.

(a) Structuration des données en dictionnaire

Le choix du stockage des annotations dans un dictionnaire en Python repose sur sa structure de données, qui stocke des paires clé-valeur. Chaque clé est unique et les valeurs peuvent être de n'importe quel type.

Dans ce cas, chaque clé du dictionnaire **annotation** est une chaîne de caractères représentant un identifiant unique pour chaque énoncé du corpus, par exemple "eslo 1", "eslo 2", etc.

La valeur associée à chaque clé est un autre dictionnaire contenant des informations sur l'énoncé correspondant. Ce dictionnaire interne a plusieurs clés telles que **source**, **nv_etudes**, **tokens**, **pos**, **neg_comp**, **y_absent**, et **schwa_absent**. Chaque clé a une liste comme valeur. Les listes **y_absent** et **schwa_absent** contiennent des booléens (**True**, **False**) ou **None**, indiquant respectivement si le phonème /y/ ou le schwa sont absents dans l'énoncé, ou si l'information n'est pas disponible.

En plus de cela, nous avons d'autres clés : **source** qui représente l'identifiant du locuteur, **nv_etudes** pour son niveau d'études, ainsi que les **tokens** et les **pos** de chaque token de l'énoncé.

Voici un exemple du format du dictionnaire :

```
annotation =
{
  "eslo 1": {
    "source": "ch_OB1",
    "nv_etudes": "bac_5_et_plus",
    "tokens": ["alors"],
    "pos": ["ADV"],
    "neg_comp": [],
    "y_absent": [None],
    "schwa_absent": [None]
  },
  "eslo 2": {
    "source": "BV1",
    "nv_etudes": "bac_1",
    "tokens": ["marche", "bien", "au", "moins", "votre", "truc", "euh"],
    "pos": ["VERB:Pres", "ADV", "ADP", "ADV", "DET", "NOUN", "ADV"],
    "neg_comp": [],
    "y_absent": [None, None, None, None, None, None, None],
    "schwa_absent": [None, None, None, None, None, None, None]
  }
}
```

Dans cet exemple, pour l'énoncé "eslo 1", le locuteur est "ch_OB1", son niveau d'études est "bac_5_et_plus", l'énoncé contient un seul token "alors" qui est un adverbe (ADV), il n'y a pas de négation complète, et l'information sur l'absence du phonème /y/ et du schwa n'est pas disponible (None).

Pour l'énoncé "eslo 2", le locuteur est "BV1", son niveau d'études est "bac_1", l'énoncé contient sept tokens avec leurs parties du discours correspondantes, il n'y a pas de négation complète, et l'information sur l'absence du phonème /y/ et du schwa pour chaque token n'est pas disponible (None).

(b) Création et formatage du fichier XML

Le format XML est un standard largement utilisé pour structurer, stocker et transporter des données. Il permet de représenter les données de manière hiérarchique et lisible aussi bien par les humains que par les machines.

Le fichier XML généré dans le cadre de notre projet suit une structure bien définie. Pour commencer, il est stocké dans une variable <DATA> dans laquelle chaque énoncé du corpus est représenté par un élément unique : <SMS> ou <ESLO>, selon le corpus, qui contient des attributs pour l'identifiant du locuteur (id) et son niveau d'études (niveau).

À l'intérieur de chaque élément, les tokens (mots ou caractères) de l'énoncé sont encapsulés dans des sous-éléments <w> avec leurs attributs respectifs : les parties

du discours (**pos**), les informations sur les négations complètes (**neg_complete**), l'absence du phonème /y/ (**y_absent**) et l'absence du schwa (**schwa**).

Voici un exemple du fichier XML avec les annotations :

```
<DATA>
  <SMS id="cmr-88milsms-p236" niveau="inconnu">
    <w pos="PROPN">Dors</w>
    <w pos="ADV">bien</w>
    <w pos="CCONJ">et</w>
    <w pos="ADV">n'</w>
    <w pos="VERB:Pres">oublie</w>
    <w pos="ADV" neg_complete="True">pas</w>
    <w pos="ADP" schwa="True">d'</w>
    <w pos="VERB:Inf">manger</w>
    <w pos="NUM">5</w>
    <w pos="NOUN">fruits</w>
    <w pos="CCONJ">et</w>
    <w pos="NOUN">legumes</w>
    <w pos="ADP">par</w>
  </SMS>
  <SMS id="cmr-88milsms-p136" niveau="inconnu">
    <w pos="PRON">Je</w>
    <w pos="ADV">n'</w>
    <w pos="VERB:Pres">arrive</w>
    <w pos="ADV" neg_complete="True">pas</w>
    <w pos="ADP">à</w>
  </SMS>
</DATA>
```

5 Résultat et retour d'expérience

Contribution : Florian, Shami

5.1 Observation des résultats

Les annotations que nous avons effectuées semblent relativement satisfaisantes. Nous aurions souhaiter évaluer de façon quantitative le travail effectué mais le manque de temps nous en a empêché.

D'après nos observations, les annotations faites semblent assez justes mais incomplètes. En effet, les règles que nous avons définies pour capturer les caractéristiques souhaitées sont

suffisamment sélectives pour éviter les grosses erreurs d’annotation mais certainement trop sélectives pour être exhaustives dans nos captures. De plus, l’annotation Spacy pourrait être meilleure si nous utilisions un plus gros modèle mais l’exécution de nos script étant longues sur nos machines, nous sommes restés sur le plus petit modèle.

L’extraction et l’uniformisation des données ont nécessité une réflexion en amont. Étant un groupe de quatre, nous avons fait des modifications à plusieurs reprises pour que les données extraits soient exploitables par le script d’annotation.

Néanmoins, malgré plusieurs essais et corrections des expressions régulières et du script d’annotation, nous n’avons pas réussi écarter certaines erreurs, surtout en ce qui concerne les annotation des schwa. Par exemple, le token « enchainement » est annoté alors qu’il ne correspond pas du tout aux conditions de captures d’absence de schwa définies dans la section 4.

5.2 Sorties

Les démarches réalisées pour extraire et modéliser les information nous ont permis dans un premier temps d’organiser et de présenter des informations que nous avons estimées pertinentes, sous des divers formats. Ainsi, nous avons en sortie cinq corpus :

- eslo_annotation.txt
- eslo_annotation.xml
- SMS_annotation.xml
- ESLO_annotation.txm
- SMS_annotation.txm

Ces données peuvent être exploitées afin d’effectuer des études statistiques des phénomènes sociolinguistiques plus avancées. Pour en citer un exemple, les corpus xml peuvent servir à faire des calculs de corrélation entre le niveau d’étude et les types de termes utilisés : des personnes avec un certain niveau d’étude sont plus ou moins susceptibles d’utiliser le schwa.

6 Conclusion

Contribution : Tout le monde

En conclusion, ce projet de recherche nous a permis d’étudier la variation sociolinguistique diastatique dans deux types de données : des données issues de l’oral et de la communication médiée par les réseaux. Les deux corpus ont été annotés en utilisant des outils tels que Spacy et des expressions régulières et les résultats obtenus ont été formalisés dans de nouvelles extensions : XML et TXM.

Ce projet nous a permis d’en apprendre plus sur ce qu’est de travailler sur la chaîne

complète de travail sur l'ingénierie des connaissances : récupération des connaissances, enrichissement / ajout de nouvelles connaissances et leur transmission. Nous avons également pu constater que ce processus demande de la réflexion et de la flexibilité. Par exemple, trouver une manière d'identifier les schwas n'a pas été de tout repos et le résultat n'est toujours pas parfait.

Surtout , cela nous a permis d'entrer de plein-pied dans le domaine du travail d'ingénieur et de réaliser notre capacité à naviguer entre la linguistique et l'informatique, en réalisant ce projet en un temps limité, tout en nous adaptant à diverses situations.

Il est vrai qu'avec plus de temps et de ressources nous aurions pu enrichir nos observations sur les résultats extraits, en faisant une analyse liée au niveau d'études, par exemple. Néanmoins, nous réussissons à aboutir à des résultats plutôt satisfaisants.

A Scripts Python utilisés

Les scripts pour l'annotation des corpus (A.1, A.2) sont identiques au niveau de l'annotation. La différence entre les deux est que dans le script d'annotation des SMS, il y a une partie pour écrire la sortie au format XML alors que pour ESLO, il y a un script à part.

A.1 Script pour l'annotation du corpus SMS et faire la sortie XML

Nom : annotation-complete.py

```

1 import sys
2 import spacy
3 import regex
4 import xml.etree.ElementTree as ET
5
6 # Vérifie que le nom du corpus à traiter est bien fourni en argument
7 if len(sys.argv) < 2:
8     print(
9         "Veuillez fournir le corpus que vous souhaitez traiter: eslo ou sms \
10         \nExemple d'exécution : python3 annotation.py eslo\n"
11     )
12     sys.exit(1)
13
14 # Récupère le nom du corpus à traiter
15 corpus = sys.argv[1]
16
17 # Ouvre les fichiers contenant les textes et les identifiants et les niveaux
18   ↪ d'étude des locuteurs
19 if corpus == "eslo":
20     texts = open("../data/transformes/xml-ESLO_contenu/concat_contenu.txt",
21                 ↪ "r")
22     ids = open("../data/transformes/xml-ESLO_id/concat_id.txt", "r")
23     niv = open("../data/transformes/tsv/eslo_id_etudes.tsv", "r")
24 elif corpus == "sms":
25     texts = open("../data/transformes/xml-SMS_contenu/SMS_contenu.txt", "r")
26     ids = open("../data/transformes/xml-SMS_id/SMS_ids.txt", "r")
27     niv = open("../data/transformes/tsv/SMS_id_education.tsv", "r")
28 else:
29     print("Corpus non reconnu")
30     sys.exit(1)
31
32 # Charge le modèle spacy pour le français
33 nlp = spacy.load("fr_core_news_md")
34 print("Modèle spacy chargé...")
35

```

```

34 # Initialise le dictionnaire qui va contenir les niveaux d'étude
35 source_niveau = {}
36
37 for line in niv:
38     source, niveau = line.split("\t")
39     source_niveau[source] = niveau.rstrip()
40
41 # reinitialise the file pointer
42 niv.seek(0)
43
44 # Initialise le dictionnaire qui va contenir les annotations
45 annotation = {}
46
47 # Liste des formes de négation et des adverbes de négation
48 list_neg = ["ne", "n'", "n"]
49 list_adv_neg = [
50     "pas",
51     "rien",
52     "jamais",
53     "point",
54     "aucunement",
55     "aucun",
56     "aucunement",
57     "nul",
58     "nullement",
59     "nulle",
60     "plus",
61 ]
62 list_neg_tok = list_neg + list_adv_neg
63
64 # Boucle sur les textes et les identifiants des locuteurs
65 n = 1
66 for t, i in zip(texts, ids):
67     # Initialise la variable qui va contenir l'information sur la négation
68     ↪ complète
69     complete_neg = None
70
71     # Affiche un message tous les 1000 textes traités
72     if n % 5000 == 0:
73         #break
74         print(f"working on eslo {n}...")
75     elif n % 1000 == 0:
76         #print(f"working on eslo {n}...")
77
78     # Ajoute les informations sur le texte
79     annotation[f"eslo {n}"] = {}
80     annotation[f"eslo {n}"][f"source"] = i.strip()

```

```

80
81     # Récupère le niveau correspondant à la source
82     niveau = source_niveau.get(i.strip())
83     if niveau:
84         annotation[f"eslo {n}"][f"niveau"] = niveau
85     else:
86         annotation[f"eslo {n}"][f"niveau"] = "inconnu"
87
88     # Initialisation des annotations
89     annotation[f"eslo {n}"]["tokens"] = []
90     annotation[f"eslo {n}"]["pos"] = []
91     annotation[f"eslo {n}"]["neg_comp"] = []
92     annotation[f"eslo {n}"]["y_absent"] = []
93     annotation[f"eslo {n}"]["schwa_absent"] = []
94
95     # Vérifie si le texte contient une négation complète
96     complete_negation_in_sent = None
97     if
98         ↪ regex.findall(r"\b(ne|n')\b.*\b({0})\b".format("|".join(list_adv_neg))),
99         ↪ t):
100         complete_negation_in_sent = True
101
102     # Traite le texte avec spacy et boucle sur les tokens
103     doc = nlp(t)
104     for j, token in enumerate(doc):
105         # Ignore les tokens vides
106         if token.pos_ != "SPACE":
107             # Ajoute le token et sa catégorie grammaticale au dictionnaire
108             ↪ d'annotations
109             annotation[f"eslo {n}"]["tokens"].append(token.text)
110             complete_pos = token.pos_
111             if token.pos_ == "VERB":
112                 if len(token.morph.get("Tense")) > 0:
113                     complete_pos = token.pos_ + ":" +
114                     ↪ token.morph.get("Tense")[0]
115                 else:
116                     complete_pos = token.pos_ + ':' + 'Inf'
117             annotation[f"eslo {n}"]["pos"].append(complete_pos)
118
119     # Ajoute l'information sur la négation complète au dictionnaire
120     ↪ d'annotations
121     if complete_negation_in_sent:
122         if token.text in list_neg:
123             position = len(annotation[f"eslo {n}"]["tokens"]) - 1
124             complete_neg = True
125
126     if (

```



```

122         complete_neg
123         and token.text in list_adv_neg
124         and regex.match(
125             r"(\^VERB.*|AUX)", annotation[f"eslo
               ↪ {n}"]["pos"][position + 1]
126         )
127     ):
128         annotation[f"eslo {n}"]["neg_comp"].append(True)
129     else:
130         annotation[f"eslo {n}"]["neg_comp"].append(None)
131
132     assert token.text == doc[j].text
133     # Ajoute l'information sur l'absence de /y/ dans "tu" au
       ↪ dictionnaire d'annotations
134     if (
135         j < len(doc) - 1
136         and token.dep_ == "nsubj"
137         and regex.match(r"t'", token.text)
138         and regex.match(r"[aeiouy]", doc[j + 1].text)
139     ):
140         if token.text == "tu":
141             annotation[f"eslo {n}"]["y_absent"].append(False)
142         else:
143             annotation[f"eslo {n}"]["y_absent"].append(True)
144
145     # Ajoute l'information sur l'absence de schwa dans les clitiques
       ↪ au dictionnaire d'annotations
146     elif (
147         j < len(doc) - 1
148         and regex.match(r"[cdjlmnqst]'", token.text)
149         and regex.match(r"^[^aeiouyââêêôôhAEIOUYÊÂÊÔÈH]", doc[j +
               ↪ 1].text)
150     ):
151         annotation[f"eslo {n}"]["schwa_absent"].append(True)
152
153     elif (
154         j < len(doc) - 1
155         and regex.match(r"^[^aeiouyââêêôôhAEIOUYÊÂÊÔÈH\']'",
               ↪ token.text)
156     ):
157         annotation[f"eslo {n}"]["schwa_absent"].append(True)
158     else:
159         annotation[f"eslo {n}"]["y_absent"].append(None)
160         annotation[f"eslo {n}"]["schwa_absent"].append(None)
161     n += 1
162
163     # print(annotation)

```

```

164
165 texts.close()
166 ids.close()
167
168 # Generate XML
169 root = ET.Element("DATA")
170
171 for key, value in annotation.items():
172     sms = ET.Element(
173         corpus.upper(), id=value["source"], niveau=value["niveau"]
174     ) # Utilisation de la variable corpus pour le nom de l'élément
175     has_content = False
176     for token, pos, neg, y, schwa in zip(
177         value["tokens"],
178         value["pos"],
179         value["neg_comp"],
180         value["y_absent"],
181         value["schwa_absent"],
182     ):
183         attributes = {"pos": pos}
184         if neg is not None:
185             attributes["neg"] = str(neg)
186         if y is not None:
187             attributes["y"] = str(y)
188         if schwa is not None:
189             attributes["schwa"] = str(schwa)
190         ET.SubElement(sms, "w", **attributes).text = token
191         has_content = True
192     if has_content:
193         root.append(sms)
194
195
196 # Prettify the XML
197 def indent(elem, level=0):
198     i = "\n" + level * "    "
199     if len(elem):
200         if not elem.text or not elem.text.strip():
201             elem.text = i + "    "
202         if not elem.tail or not elem.tail.strip():
203             elem.tail = i
204         for elem in elem:
205             indent(elem, level + 1)
206         if not elem.tail or not elem.tail.strip():
207             elem.tail = i
208     else:
209         if level and (not elem.tail or not elem.tail.strip()):
210             elem.tail = i

```

```

211
212
213 indent(root)
214
215 # Write the pretty XML to a file
216 output_file = f"../data/annotate/{corpus.upper()}_annotation.xml" #
    ↪ Utilisation de la variable corpus pour le nom du fichier de sortie
217 tree = ET.ElementTree(root)
218 with open(output_file, "wb") as fh:
219     tree.write(fh, encoding="utf-8", xml_declaration=True)
220
221 print(
222     f"Fichier XML joliment formaté pour le corpus {corpus.upper()} généré
    ↪ avec succès."
223 )
224

```

A.2 Script pour l'annotation du corpus ESLO

Nom : annotation.py

```

1 import sys
2 import spacy
3 import regex
4
5 # Vérifie que le nom du corpus à traiter est bien fourni en argument
6 if len(sys.argv) < 2:
7     print(
8         "Veuillez fournir le corpus que vous souhaitez traiter: eslo ou sms \
9         \nExemple d'exécution : python3 annotation.py eslo\n"
10    )
11    sys.exit(1)
12
13 # Récupère le nom du corpus à traiter
14 corpus = sys.argv[1]
15
16 # Ouvre les fichiers contenant les textes et les identifiants et les niveaux
    ↪ d'étude des locuteurs
17 if corpus == "eslo":
18     texts = open("../data/transformes/xml-ESLO_contenu/concat_contenu.txt",
    ↪ "r")
19     ids = open("../data/transformes/xml-ESLO_id/concat_id.txt", "r")
20     nv_etudes = open("../data/transformes/xml-eslo_nvetudes.txt", "r")
21 elif corpus == "sms":
22     texts = open("../data/transformes/xml-SMS_contenu/SMS_contenu.txt", "r")

```

```
23     ids = open("../data/transformes/xml-SMS_id/SMS_ids.txt", "r")
24 else:
25     print("Corpus non reconnu")
26     sys.exit(1)
27
28 # Charge le modèle spacy pour le français
29 nlp = spacy.load("fr_core_news_md")
30 # print("Modèle spacy chargé...")
31
32 # Initialise le dictionnaire qui va contenir les annotations
33 annotation = {}
34
35 # Liste des formes de négation et des adverbes de négation
36 list_neg = ["ne", "n'", "n"]
37 list_adv_neg = [
38     "pas",
39     "rien",
40     "jamais",
41     "point",
42     "aucunement",
43     "aucun",
44     "aucunement",
45     "nul",
46     "nullement",
47     "nulle",
48     "plus",
49 ]
50 list_neg_tok = list_neg + list_adv_neg
51
52 # Boucle sur les textes et les identifiants des locuteurs
53 n = 1
54 for t, i, ne in zip(texts, ids, nv_etudes):
55     # Initialise la variable qui va contenir l'information sur la négation
56     ↪ complète
57     complete_neg = None
58
59     # if n % 1000 == 0:
60     #     break
61     # elif n % 100 == 0:
62     #     print(f"working on eslo {n}...")
63
64     # Ajoute les informations sur le texte et initialisation des annotations
65     annotation[f"eslo {n}"] = {}
66     annotation[f"eslo {n}"]["source"] = i.strip()
67     annotation[f"eslo {n}"]["nv_etudes"] = ne.strip()
68     annotation[f"eslo {n}"]["tokens"] = []
69     annotation[f"eslo {n}"]["pos"] = []
```

```

69     annotation[f"eslo {n}"]["neg_comp"] = []
70     annotation[f"eslo {n}"]["y_absent"] = []
71     annotation[f"eslo {n}"]["schwa_absent"] = []
72     # print("\t<eslo>")
73
74     # Vérifie si le texte contient une négation complète
75     complete_negation_in_sent = None
76     if
77         ↪ regex.findall(r"\b(ne|n')\b.*\b({0})\b".format("|".join(list_adv_neg)),
78         ↪ t):
79         complete_negation_in_sent = True
80
81     # Traite le texte avec spacy et boucle sur les tokens
82     doc = nlp(t)
83     for j, token in enumerate(doc):
84         # Ignore les tokens vides
85         if token.pos_ != "SPACE":
86             # Ajoute le token et sa catégorie grammaticale au dictionnaire
87             ↪ d'annotations
88             annotation[f"eslo {n}"]["tokens"].append(token.text)
89             # print(f'<w pos="{token.pos_}", end=''')
90             complete_pos = token.pos_
91             if token.pos_ == "VERB":
92                 if len(token.morph.get("Tense")) > 0:
93                     complete_pos = token.pos_ + ":" +
94                     ↪ token.morph.get("Tense")[0]
95                 else:
96                     complete_pos = token.pos_ + ":" + "Inf"
97             annotation[f"eslo {n}"]["pos"].append(complete_pos)
98
99     # Ajoute l'information sur la négation complète au dictionnaire
100     ↪ d'annotations
101     if complete_negation_in_sent:
102         if token.text in list_neg:
103             position = len(annotation[f"eslo {n}"]["tokens"]) - 1
104             complete_neg = True
105
106         if (
107             complete_neg
108             and token.text in list_adv_neg
109             and regex.match(
110                 r"(\bVERB.*|AUX)", annotation[f"eslo
111                 ↪ {n}"]["pos"][position + 1]
112             )
113         ):
114             annotation[f"eslo {n}"]["neg_comp"].append(True)
115         else:

```

```

110         annotation[f"eslo {n}"]["neg_comp"].append(None)
111
112     assert token.text == doc[j].text
113     # Ajoute l'information sur l'absence de /y/ dans "tu" au
114     ↪ dictionnaire d'annotations
115     if (
116         j < len(doc) - 1
117         and token.dep_ == "nsubj"
118         and regex.match(r"t'", token.text)
119         and regex.match(r"[aeiouy]", doc[j + 1].text)
120     ):
121         if token.text == "tu":
122             annotation[f"eslo {n}"]["y_absent"].append(False)
123             print(token, doc[j + 1].text, " : ypresent")
124         else:
125             annotation[f"eslo {n}"]["y_absent"].append(True)
126             print(token, doc[j + 1].text, " : yabsent")
127
128     # Ajoute l'information sur l'absence de schwa dans les clitiques
129     ↪ au dictionnaire d'annotations
130     elif (
131         j < len(doc) - 1
132         and regex.match(r"[cdjlmnqst]'", token.text)
133         and regex.match(r"^[^aeiouyââéèêôhAEIOUYÉÂÂÊÔÈH]", doc[j +
134         ↪ 1].text)
135     ):
136         annotation[f"eslo {n}"]["schwa_absent"].append(True)
137     elif (
138         j < len(doc) - 1
139         and regex.match(r"^j[^aeiouyââéèêôhAEIOUYÉÂÂÊÔÈH']'",
140         ↪ token.text)
141     ):
142         annotation[f"eslo {n}"]["schwa_absent"].append(True)
143     else:
144         annotation[f"eslo {n}"]["y_absent"].append(None)
145         annotation[f"eslo {n}"]["schwa_absent"].append(None)
146
147     n += 1
148
149 print(annotation)
150
151 texts.close()
152 ids.close()
153

```

A.3 Script pour l'extraction du corpus sms

Nom : sms_extraction.py

```

1 from collections import defaultdict
2 from lxml import etree
3 from pprint import pprint
4 import regex
5 import os
6
7 ## suppression des espaces multiples
8 def clean(text):
9     text = regex.sub(r'\s+', ' ', text.strip(), flags=regex.MULTILINE)
10    return text
11
12 ## recuperation du root
13 def get_root():
14     root =
15     ↪ etree.parse('../data/brut/xml-SMS/cmr-88milsms-tei-v1.xml').getroot()
16    return root
17
18
19
20 def get_person_niveauEtude():
21     """
22     Recupere le niveau d'education de chaque personne
23     :return: dictionnaire {id : niveau_education}
24     : exporte dans un fichier tsv
25     """
26     root = get_root()
27     person_edu = {}
28     for p, edu in zip(root.xpath('//*[name()="person"]'),
29     ↪ root.xpath('//*[name()="education"]')):
30         xml_id = p.get('{http://www.w3.org/XML/1998/namespace}id')
31         education = edu.text
32         person_edu[xml_id] = education
33
34     # export to tsv file xml_id \t education
35     # create filepath if not exists
36     if not os.path.exists('../data/transformes/tsv/'):
37         os.makedirs('../data/transformes/tsv/')
38
39     ''' export to tsv file'''
40     ## décommenter si on souhaite exporter à un fichier tsv
41     with open('../data/transformes/tsv/SMS_id_education.tsv', 'a') as f:
42         f.write(f"{xml_id}\t{education}\n")

```

```

42
43     return person_edu
44
45
46 def append_to_tsv(who, edu):
47     """
48     AJOUTE les locuteurs avec niveau d'études inconnus au fichier tsv
49     :param who: id de la personne
50     :param edu: niveau d'education
51     """
52     with open('../data/transformes/tsv/SMS_id_education.tsv', 'a') as f:
53         f.write(f"{who}\t{edu}\t\n")
54
55 def create_id_dialogues_txt(root, person_edu):
56     """
57     Creation des fichiers txt pour id et contenu en parallele
58     :param root: root du fichier xml
59     :param person_edu: dictionnaire {id : niveau_education}
60
61     :return: dictionnaire { who : [dialogues] }
62     : stockage dans des fichiers tsv
63
64     """
65     dialogues = defaultdict(lambda: defaultdict(list))
66     without_edu_info = []
67     # recup { who : [dialogues] }
68     for post in root.xpath('//*[name()="post"]'):
69         who = post.get('who')[1:]
70         for p in post.xpath('..//following-sibling::*[name()="p"]'):
71
72             ## création fichier txt pour id et contenu
73             try:
74                 # create a txt file to store the id
75                 if len(clean(p.text)) != 0: # exporter seulement si non vides
76                     if not os.path.exists('../data/transformes/xml-SMS_id/'):
77                         os.makedirs('../data/transformes/xml-SMS_id/')
78                     with
79                         ↪ open('../data/transformes/xml-SMS_id/SMS_ids.txt',
80                         ↪ 'a') as f:
81                         f.write(f"{who}\n")
82
83             if not
84                 ↪ os.path.exists('../data/transformes/xml-SMS_contenu/'):
85                 os.makedirs('../data/transformes/xml-SMS_contenu/')
86             with
87                 ↪ open('../data/transformes/xml-SMS_contenu/SMS_contenu.txt',
88                 ↪ 'a') as f:
89                 f.write(f"{clean(p.text)}\n")

```



```

84
85         # creation du dictionnaire {id : dialogue}
86         '''
87         if who not in dialogues.keys():
88             dialogues[who] = [clean(p.text)]
89         else:
90             dialogues[who].append(clean(p.text))
91         '''
92     except Exception as e:
93         print(f"An error occurred while updating dialogues for
94         ↪ '{who}': {e}")
95
96     ## creation des fichiers tsv pour id education et dialogues
97     try:
98         if len(clean(p.text)) != 0: # exporter seulement si non vides
99             with
100                 ↪ open('../data/transformes/tsv/SMS_id_education_dialogues.tsv',
101                 ↪ 'a') as f:
102                 if who not in person_edu.keys():
103                     person_edu[who] = 'inconnu'
104                     append_to_tsv(who, person_edu[who])
105
106                 ↪ f.write(f"{who}\t{person_edu[who]}\t{clean(p.text)}\n")
107
108     except Exception as e:
109         # print(f"An error occurred while updating dialogues for
110         ↪ '{who}': {e}")
111         without_edu_info.append(who)
112
113     return dialogues
114
115 def main():
116     root = get_root() # recuperation du root
117     person_edu = get_person_niveauEtude() # dictionnaire {id :
118         ↪ niveau_education}
119     dialogues = create_id_dialogues_txt(root, person_edu) # creation des
120         ↪ fichiers txt pour id et contenu
121
122 if __name__ == '__main__':
123     main()
124

```

A.4 Script pour l'extraction du corpus ESLO

Nom : `extract_eslo.py`

```

1  # Ce script sert à extraire les ids et le contenu du corpus ESLO
2  import os
3  from xml.etree import ElementTree as ET
4
5  # Chemin du dossier contenant les fichiers XML
6  input_folder = '../data/xml-ESLO_nettoyee/'
7  # Dossiers de sortie pour les fichiers d'ids et de contenus
8  output_folder_who = '../data/xml-ESLO_id/'
9  output_folder_content = '../data/xml-ESLO_contenu/'
10
11 # Vérifie si les dossiers de sortie existent, sinon les crée
12 os.makedirs(output_folder_who, exist_ok=True)
13 os.makedirs(output_folder_content, exist_ok=True)
14
15 # Parcours des fichiers dans le dossier d'entrée
16 for filename in os.listdir(input_folder):
17     if filename.endswith('.xml'):
18         # Chemin complet du fichier XML d'entrée
19         input_file_path = os.path.join(input_folder, filename)
20         tree = ET.parse(input_file_path)
21         root = tree.getroot()
22
23         # Chemins de sortie pour les fichiers who_values.txt et
24         ↪ content_values.txt
25         who_output_path = os.path.join(output_folder_who,
26         ↪ filename.replace('.xml', '_who.txt'))
27         content_output_path = os.path.join(output_folder_content,
28         ↪ filename.replace('.xml', '_contenu.txt'))
29
30         # Ouvre les fichiers de sortie
31         with open(who_output_path, 'w') as who_file,
32         ↪ open(content_output_path, 'w') as content_file:
33             # Parcours des enfants de la racine
34             for child in root:
35                 # Vérifie si l'élément est une balise <u>
36                 if child.tag == 'u':
37                     # Récupère la valeur de l'attribut "who"
38                     who_value = child.attrib.get('who', '')
39                     # Récupère le contenu de la balise <u>, vérifie si c'est
40                     ↪ None avant d'appeler strip()
41                     content_value = child.text.strip() if child.text is not
42                     ↪ None else ''
43                     # Écrit les valeurs dans les fichiers de sortie

```

```

38         who_file.write(who_value + '\n')
39         content_file.write(content_value + '\n')
40
41     print("Traitement terminé.")
42

```

A.5 Script pour la création d'un fichier niveau d'études pour ESLO

Nom : add_etudes_eslo.py

```

1  # Ce script a servi à faire un fichier correspondant aux niveaux d'études des
   ↳ locuteurs ESLO
2  import csv
3
4  # Fichiers d'entrée
5  id_file = '../data/transformes/xml-ESLO_id/concat_id.txt'
6  level_file = '../data/transformes/tsv/eslo_id_etudes.csv'
7  output_file = '../eslo_nvetudes.txt'
8
9  # Lecture du fichier des niveaux d'études
10 levels = {}
11 with open(level_file, mode='r', encoding='utf-8') as csvfile:
12     csvreader = csv.reader(csvfile, delimiter=',')
13     for row in csvreader:
14         if len(row) >= 2:
15             levels[row[0]] = row[1] # ID est la première colonne, niveau
                                   ↳ d'études est la deuxième colonne
16
17 # Lecture du fichier des IDs et ajout des niveaux d'études
18 with open(id_file, mode='r', encoding='utf-8') as infile, open(output_file,
   ↳ mode='w', encoding='utf-8') as outfile:
19     for line in infile:
20         id_ = line.strip()
21         level = levels.get(id_, 'Unknown') # En cas d'ID non trouvé dans le
                                   ↳ fichier CSV
22         outfile.write(f"{level}\n")
23
24 print("Traitement terminé. Vérifiez le fichier 'eslo_nvetudes.txt'.")

```

A.6 Script pour transformer le fichier ESLO annoté en format XML

Nom : eslo_to_xml.py

```
1 # Ce script sert à transformer le corpus ESLO annoté en format XML, avec de
  ↪ jolies balises.
2
3 import ast
4 import xml.etree.ElementTree as ET
5 import xml.dom.minidom
6
7 # Lire le contenu du fichier TXT
8 with open('../data/annotate/eslo_annotation.txt', 'r', encoding='utf-8') as
  ↪ file:
9     data = file.read()
10
11 # Convertir le contenu du fichier en dictionnaire Python
12 data_dict = ast.literal_eval(data)
13
14 # Créer l'élément racine
15 root = ET.Element('data')
16
17 # Traiter chaque entrée dans le dictionnaire
18 for key, value in data_dict.items():
19     eslo_id = value.get('source', 'unknown')
20     nv_etudes = value.get('nv_etudes', 'unknown')
21
22     # Créer un élément 'eslo' avec les attributs 'id' et 'niveau'
23     eslo_element = ET.SubElement(root, 'eslo', id=eslo_id, niveau=nv_etudes)
24
25     # Ajouter les tokens et leurs pos en tant qu'éléments 'w'
26     tokens = value.get('tokens', [])
27     pos = value.get('pos', [])
28     neg_comp = value.get('neg_comp', [])
29     y_absent = value.get('y_absent', [])
30
31     for i in range(len(tokens)):
32         token = tokens[i]
33         pos_tag = pos[i]
34
35         # Créer l'élément 'w' avec l'attribut 'pos'
36         w_element = ET.SubElement(eslo_element, 'w', pos=pos_tag)
37         w_element.text = token
38
39         # Ajouter les attributs neg_comp et y_absent si présents
40         if i < len(neg_comp) and neg_comp[i] is not None:
41             w_element.set('neg_comp', str(neg_comp[i]))
42         if i < len(y_absent) and y_absent[i] is not None:
43             w_element.set('y_absent', str(y_absent[i]))
44
45 xml_str = ET.tostring(root, encoding='unicode')
```

```
46
47 # Utiliser minidom pour prettifier le XML
48 dom = xml.dom.minidom.parseString(xml_str)
49 pretty_xml_str = dom.toprettyxml(indent="    ")
50
51 # Écrire le contenu XML dans un fichier avec une mise en page correcte
52 with open('../data/annotate/eslo_annotation.xml', 'w', encoding='utf-8') as
    ↪ file:
53     file.write(pretty_xml_str)
54
55 print("Annotation XML terminée.")
56
```