

Remerciement

Je tiens à travers ce rapport a exprimé mes sincré remerciement a tout l'équipe pédagogique de ESTF , Je tiens à exprimer ma profonde reconnaissance envers les professionnels responsables de la filière Génie informatique pour m'avoir offert une formation de qualité.

Je souhaite également exprimer ma gratitude envers mon encadrant de PFE, le professeur **M. BENSLIMANE Mohammed**, pour ses précieux conseils, sa disponibilité et son soutien tout au long de ce projet, et qui n'a ménagé aucun effort pour m'aider et me soutenir, ainsi que pour ses conseils précieux et ses directives pertinentes.

Je tiens aussi à remercier les membres du jury qui ont pris le temps de participer à la soutenance et d'évaluer mon travail. Enfin, je souhaite exprimer ma gratitude à tous ceux qui ont contribué, de près ou de loin, à la réussite de ce travail. Un grand merci du fond du cœur.

Table des matières

I. Cahier de charge	3
Introduction	3
1. Présentation du projet	4
2. Objectifs du projet	4
2.1. À qui s'adresse ce projet	6
2.2. Etat actuel des ressources et fonctionnement du préexistant	7
2.3. Conception	11
3. Présentation attendue	11
3.1. Création, réalisation et récupération de contenu	11
3.2. Mis à jour	18
4. Présentation et organisation du travail	18
4.1. Planning prévisionnel	18
5. Livrables attendue	19
Conclusion	25
II. Conception du projet	25
III. Outils et technologie	26
IV. Réalisation	29
V. Conclusion	34

I. Cahier des charges

Introduction

Ces dernières années ont vu une montée en importance de deux sous-domaines fondamentaux de l'intelligence artificielle : l'apprentissage automatique et le traitement automatique du langage naturel (NLP). L'apprentissage automatique révolutionne la manière dont les systèmes informatiques peuvent apprendre à partir de données et améliorer leurs performances au fil du temps sans programmation explicite. Grâce à l'utilisation de modèles statistiques et d'algorithmes, l'apprentissage automatique permet aux ordinateurs d'identifier des schémas complexes dans les données, leur permettant de prendre des décisions autonomes et d'accomplir efficacement des tâches spécifiques.

Au même sens, le traitement automatique du langage, ou NLP, se concentre sur la manière dont les machines peuvent comprendre et interagir naturellement avec le langage humain. Des techniques avancées telles que l'analyse syntaxique, la reconnaissance vocale et la traduction automatique permettent au traitement automatique du langage naturel (NLP) de permettre aux ordinateurs de comprendre et de produire du langage humain dans une variété de contextes. Cela ouvre la porte à une large gamme d'applications, notamment les chatbots, les systèmes de traduction automatique et les systèmes de recommandation de contenu.

En combinant le NLP & le machine_learning, nous créons une collaboration puissante qui permet aux machines de comprendre, d'interpréter et de produire du langage naturel de manière plus avancée, tel que l'amélioration de la capacité des machines à traiter saisir le sens et le contexte du langage humain, ainsi qu'à générer des réponses pertinentes et naturelles. Cela ouvre de nouvelles perspectives passionnantes pour l'avenir de l'intelligence artificielle et de ses applications dans divers domaines.

1. Présentation du projet

Avec la multiplication des contenus disponibles sur Internet, il est devenu essentiel de comprendre et d'analyser rapidement les informations textuelles. Dans un environnement en constante évolution, mon projet s'est concentré sur le développement d'une solution complète et efficace, exploitant les techniques avancées de traitement automatique du langage naturel (NLP).

L'objectif principal de mon projet était de fusionner les technologies de NLP et d'apprentissage automatique pour créer une application capable d'identifier automatiquement la langue d'un contenu audio, puis de le transcrire en texte pour une analyse ultérieure afin de donner un résumé complet., et cela avec une intégration des algorithmes de NLP et des modèles d'apprentissage automatique pour accomplir cette tâche et découvrir des caractéristiques avancées permettant de reconnaître et de traiter les différentes langues parlées dans les fichiers audio, puis de les transcrire avec une grande précision.

Une fois le contenu audio converti en texte, l'application doit produire un résumé concis et informatif du contenu extrait, offrant ainsi aux utilisateurs une vue d'ensemble rapide et pertinente des informations contenues dans le fichier audio.

2. Objectifs du projet

L'objectif principal de mon projet était de fusionner les technologies de NLP et d'apprentissage automatique pour créer une application capable d'identifier automatiquement la langue d'un contenu audio, puis de le transcrire en texte pour une analyse ultérieure afin de donner un résumé complet.

il peut s'apparaître que le projet vise à réaliser deux tâches distinctes, mais elles sont également interconnectées, par ce que en tant que utilisateur international, qui cherche d'un résumé afin d'obtenir une information exacte importante, d'après

un document qui peut être dans une langue incompréhensible pour lui, c'est pour cela ,La détection de langue permet d'obtenir une expérience personnalisée en recevant du contenu dans leur langue préférée et facilite la classification et l'organisation des données dans des environnements multilingues.

La Fonctionnalité de détection de langue revêt une importance capitale dans divers contextes, notamment pour les services de traduction automatique, les moteurs de recherche multilingues et les plateformes de médias sociaux. La détection de langue permet aux utilisateurs de bénéficier d'une expérience personnalisée en recevant du contenu dans leur langue préférée, et facilite également la classification et l'organisation des données dans des environnements multilingues.

D'autre part , La fonctionnalité de résumé automatique est particulièrement utile dans les cas où les utilisateurs ont besoin d'obtenir rapidement une compréhension générale du contenu sans avoir à lire ou écouter l'intégralité du texte ou de l'audio.

Cette combinaison présente de nombreux bénéfices, tels qu'une expérience utilisateur optimisée, une productivité renforcée et une compréhension améliorée du contenu dans un contexte multilingue.

Voici quelque avantage :

Personnalisation et Accessibilité : L'application peut détecter automatiquement la langue dans laquelle le contenu est rédigé et générer un résumé dans cette langue, ce qui offre une expérience personnalisée aux utilisateurs et facilite l'accès à l'information dans leur langue préférée.

Gain de Temps : En générant automatiquement des résumés, l'application permet aux utilisateurs de gagner du temps en évitant de devoir parcourir l'intégralité du contenu pour en extraire les informations pertinentes.

Analyse Multilingue : L'application peut être utilisée pour analyser et résumer des contenus provenant de différentes langues, ce qui facilite la comparaison et l'analyse de données dans un contexte multilingue.

2.1. À qui s'adresse ce projet

le projet s'adresse à diverses parties prenantes et peut bénéficier à différents utilisateurs, il peut être :

Utilisateurs Individuels : Les personnes qui souhaitent obtenir un résumé rapide d'un contenu textuel dans une langue spécifique sans avoir à lire tout le contenu. Cela peut inclure des étudiants, des professionnels pressés ou toute personne cherchant à gagner du temps lors de la consommation de contenu audio.

Entreprises de Médias et de Veille : Les entreprises qui traitent de grandes quantités de contenu et qui ont besoin de méthodes efficaces pour extraire les informations essentielles. Alors, le projet peut leur fournir des résumés automatiques de manière rapide et efficace, ce qui facilite leur analyse et leur prise de décision.

Services de Traduction et d'Interprétation : Les services de traduction automatique et d'interprétation peuvent l'utiliser pour détecter automatiquement la langue dans laquelle un contenu audio est rédigé, ce qui peut faciliter le processus de traduction en identifiant rapidement la langue cible.

Équipes de Recherche et d'Analyse : Les équipes de recherche et d'analyse qui travaillent sur des projets multilingues peuvent utiliser ce projet afin d'analyser et résumer efficacement des contenus audio dans différentes langues, ce qui leur permet d'obtenir rapidement des informations pertinentes.

Développeurs et Innovateurs : Les développeurs intéressés par le domaine du traitement automatique du langage naturel et de l'apprentissage automatique peuvent utiliser notre projet comme une base pour développer de nouvelles applications et solutions dans ce domaine.

Le projet offre une solution polyvalente et efficace pour répondre aux besoins diversifiés d'utilisateurs et d'organisations. En combinant la détection de langue, la conversion audio en texte et la génération de résumés automatiques, notre application permet d'optimiser la gestion et l'analyse de contenu multilingue. Cela offre une valeur ajoutée significative dans des domaines variés tels que la recherche d'information, l'analyse de contenu et bien d'autres encore.

2.2. Etat actuel des ressources et fonctionnement du préexistant

Avant de commencer dans le processus de développement de projet, il faut premièrement évaluer l'état actuel des ressources disponibles ainsi que le fonctionnement du système préexistant, le cas échéant. Cette évaluation nous permettra de comprendre les fondations sur lesquelles le projet sera basé dans son réalisation , et de déterminer les défis potentiels à rencontrer en cours de route.

Dans cette section, il y a une examination ressources disponibles, telles que les données d'entraînement, et les technologies disponibles.

Les ressources actuelles :

Données d'Entraînement :

Le projet est divisé en trois parties distinctes. La première concerne la détection de langue d'un contenu audio ou textuel, la seconde se concentre sur le résumé automatique, tandis que la troisième partie porte sur la génération de résumés de texte. En termes de données d'entraînement, la partie audio nécessite une quantité significative de données, suivie de la partie texte.

Cependant, la partie de résumé automatique ne nécessite pas spécifiquement de données d'entraînement.

Après des recherches approfondies, voici les résultats obtenus :

Pour la partie audio, qui requiert une importante quantité de données, j'ai découvert une base de données contenant des enregistrements audio dans trois langues spécifiques : le français, l'allemand et l'anglais. Cette base de données est la plus conséquente que j'ai trouvée dans la communauté Kaggle, avec une taille de 16 Go.

pour la partie text (detection de langue), sur le même site, j'ai trouvé un fichier CSV contenant un ensemble de phrases dans plusieurs langues, chaque phrase étant associée à sa langue respective. Les langues couvertes comprennent l'arabe, le français, l'anglais, l'allemand, le russe, l'hindi, entre autres. La taille totale du fichier est de 1 911 016 octets.

Technologies Disponibles :

la partie audio :

Tout d'abord, la transformation de l'audio en spectrogramme Mel-Fréquence Cepstral Coefficients (MFCC). Ensuite, l'application d'un réseau neuronal (NN) afin de déterminer les caractéristiques fréquentielles de chaque langue. La réalisation de cela nécessite la présence des bibliothèques suivantes :

NumPy et SciPy pour le traitement des données et les calculs scientifiques :

Ils sont deux bibliothèques essentielle pour le traitement des données qui offrent des structures de données puissantes pour manipuler des tableaux multidimensionnels ainsi que des fonctions pour effectuer des opérations

mathématiques sur ces tableaux, il sera utilisé pour représenter les données audio sous forme de tableaux et effectuer des opérations de transformation et de traitement nécessaires pour obtenir les MFCC spectrogram

Matplotlib pour la visualisation des spectrogrammes :

Cette représentation visuelle nous permet d'analyser les caractéristiques fréquentielles par laquelle nous serons capable d'identifier la langue de contenu audio.

TensorFlow avec Keras pour la création du réseau neuronal (NN) :

Ces deux bibliothèques Ensemble, offrent un environnement puissant pour la création, l'entraînement et le déploiement de modèles de ML. Et ils nous permet de créer et entraîner un réseau neuronal qui analyse les caractéristiques fréquentielles des spectrogrammes.

La partie text :

Pour traiter le contenu textuel, il faut d'abord le préparer avant d'appliquer le modèle de machine learning qui nous permettra de détecter la langue. Cette préparation implique l'application de techniques de traitement du langage naturel (NLP), suivie de la visualisation de la base de données. Pour ce faire :

Pandas : Pandas offre des structures de données flexibles et des outils pour manipuler et analyser des données. Nous l'utiliserons pour visualiser la base de données textuelles, notamment pour explorer ses caractéristiques et sa distribution.

NLTK (Natural Language Toolkit) : Nous l'utiliserons pour effectuer les traitements NLP nécessaires, tels que le nettoyage du texte, la tokenization, le stemming ou le lemmatization. Ces étapes sont essentielles pour préparer les données textuelles avant de les utiliser dans le modèle de détection de langue.

Math : La bibliothèque Math est une bibliothèque standard de Python qui fournit des fonctions mathématiques de base. Elle peut être utilisée pour certaines opérations mathématiques nécessaires dans l'implémentation du modèle choisi.

Pour le résumé automatique:

Dans ce cas, en utilisant uniquement NLTK et la bibliothèque mathématique de Python, puisque nous serons implémenter une méthode de résumé automatique qui se concentre sur le traitement du langage naturel (NLP) et appliquer les techniques mathématiques pour générer des résumés.

NLTK (Natural Language Toolkit) : pour effectuer les traitements de base du texte, tels que la tokenization, la lemmatization, la suppression des mots vides (stopwords), et éventuellement l'analyse de la syntaxe pour identifier les phrases clés ou les relations entre les mots. Ces étapes permettront de préparer le texte avant la génération du résumé.

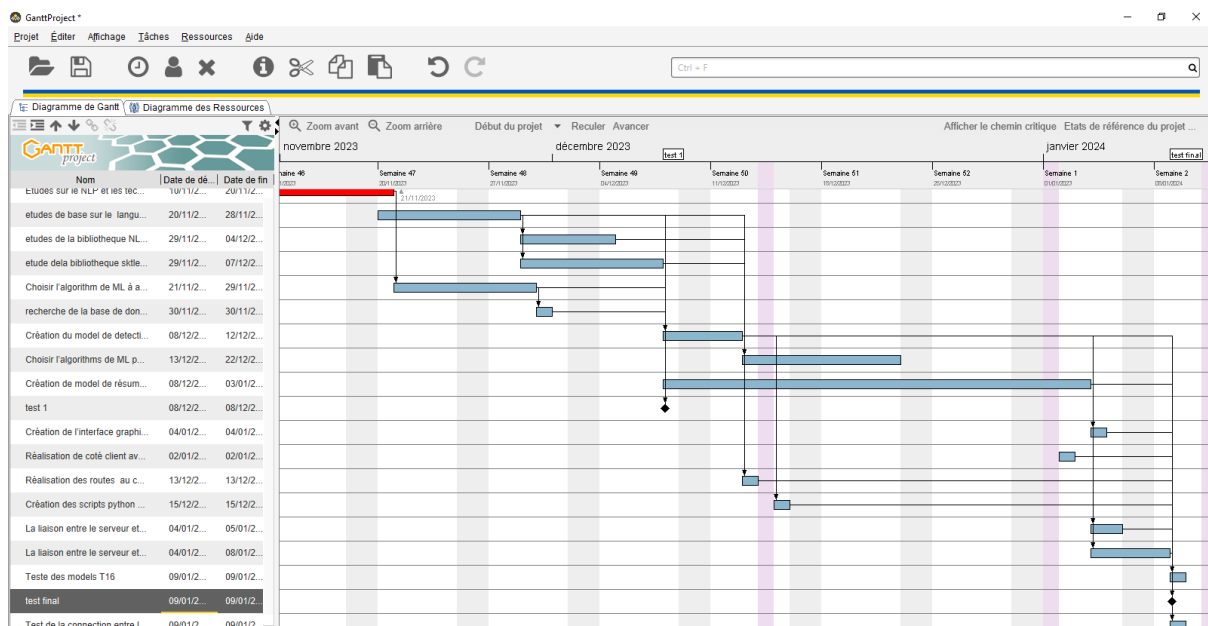
Bibliothèque mathématique (Math) : pour les techniques mathématiques qui sert à évaluer l'importance des phrases ou des mots dans le texte

Pour l'interface graphic :

Pour créer une interface utilisateur, il faut d'abord mettre en place des routes dans un serveur Node.js. Ensuite, nous utilisons AJAX pour communiquer avec le serveur côté client. Pour intégrer notre modèle Python, nous devons établir une communication entre le serveur Node.js et notre script Python. Dans mon cas, j'ai utilisé le module `child_process`, qui est intégré à Node.js et permet de créer des processus enfants. Ces processus enfants sont des processus indépendants du processus principal de Node.js et peuvent exécuter des commandes dans le shell du système d'exploitation. Nous utilisons spécifiquement la méthode `spawn` pour envoyer les paramètres nécessaires à l'exécution du modèle et récupérer le résultat pour l'afficher dans l'interface graphique. Pour la partie front-end, nous utilisons HTML et CSS pour créer une interface statique.

2.3. conception

Image représente le diagramme de gantt de projet :



3. Présentation attendues

3.1 Création, réalisation et récupération de contenu

La création, la réalisation et la récupération de contenu représentent des étapes fondamentales pour alimenter notre système de détection de langue et de résumé automatique. Dans cette section, nous détaillons les processus impliqués dans la collecte, la création ou l'acquisition de données textuelles et audio pour notre application.

Collecte de Données :

Comme déjà mentionné dans la partie 1.4 ,on a obtenu les données de training et de testing d'après la plateforme kaggle,au côté audio c'est une base de données (audio de format .wav)qui contient 16GB des audio en trois langues l'allemand ,le français et l'anglais, pour chaque langue ,mal et femal, et pour le côté textuelle ,on a utilisé une base de données qui contient des phrases en différentes langues tel que AR,FR,E,EN,SP..),chaque ligne deux colonnes une de langues et l'autre de la phrase.

Création de contenu :

Les étapes nécessaires afin de créer notre projet sont les suivantes :

Partie1 : détection de langue d'un audio

Dans cette partie on doit premièrement préparer les données afin d'appliquer un NN sur elle,la préparation des données est l'équivalence de la transformation des audio en MFCC spectrogram

Définition :

Un spectrogramme est la représentation visuelle d'un son,qui est basé sur le changement de la fréquence par le temps.

pourquoi les spectrogram ?

les signaux contiennent plusieurs composantes fréquentielles ,qui changent d'une personne à une autre,et d'une langue à une autre ,car on peut trouver que la tonalité d'une personne se change entre eux langues différentes,alors c'est la meilleure notion qu'on peut utiliser et étudier afin de classer les audio. Et pour mieux avoir les caractéristiques fréquentielles de l'audio,le spectrogramme reste le meilleur choix , les étapes de réalisation de spectrogramme :

ans cette partie nous serons voir 3 aspects, pratique , mathématique et aspect traitement (code) avec python

les formats de signal audio sont nombreux, .WAV, .FLAC, MP3, chacun utilisé dépend du besoin :

.MP3 : format audio compressé qui réduit la taille des fichiers audio en éliminant certaines données auditives moins essentielles . Il est utilisé pour la diffusion de musique en ligne et le stockage de fichiers audio sur des appareils portables puisqu'il est d'une taille réduite.

.FLAC : un format audio sans perte, qui comprime les données audio sans perte de qualité et utilisé par les audiophiles pour archiver leur musique dans une qualité audio haute fidélité.

.WAV : le format standard pour les fichiers audio sur les ordinateurs Windows. Les données audio sont non compressées, et offrent une qualité audio élevée, mais occupent plus d'espace sur le disque dur.

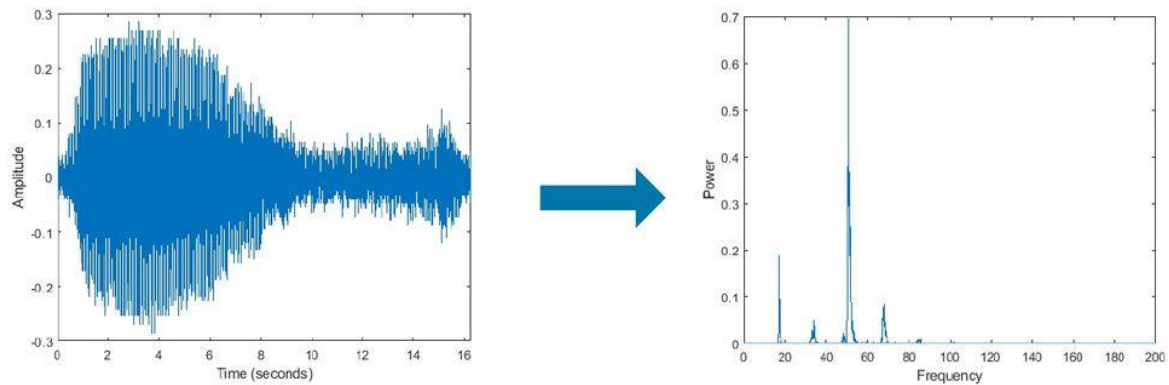
- afin d'étudier toutes les caractéristiques d'un signal audio, on utilise un format qui sert à garder les fonctionnalités de l'audio, c-à-d .WAV ou .FLAC , dans mon cas j'ai obtenu un ensemble des données qui sont des audios en trois langues (GERMAN, ENGLISH, FRENCH), dans chaque langue il y a une partie female et une partie male, c'est de 16GB , dans laquelle les audios sont en format .FLAC.

Les étapes :

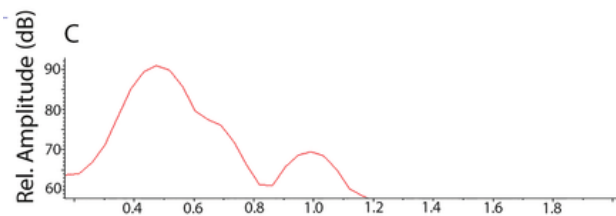
1. le fenêtrage : premièrement en commande à découper notre signal en plusieurs petites fenêtres, chaque fenêtre a une taille spécifique, et la durée entre deux frames successives est nommée hop size , lors du fenêtrage, on utilise la méthode de overlapping, c-à-d la deuxième fenêtre est incluse dans la première, autrement dit , le hop size est inférieur à la taille de la trame, on utilise cette méthode pour ne pas perdre aucune partie de notre signal audio.

2. FFT :

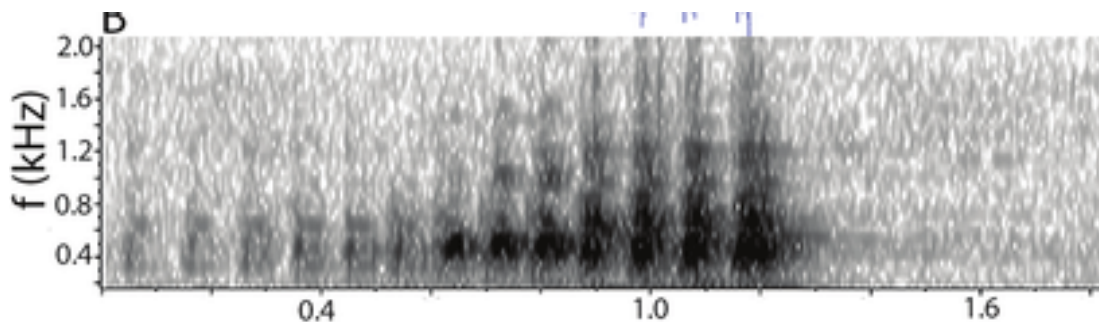
Le mot FFT est une abréviation de (Fast Fourier Transform), il nous permet de passer de la représentation de signal en fonction de temps vers sa représentation en fonction de la fréquence, c'est le même concept du spectrogramme , pour une fenêtre uniquement.



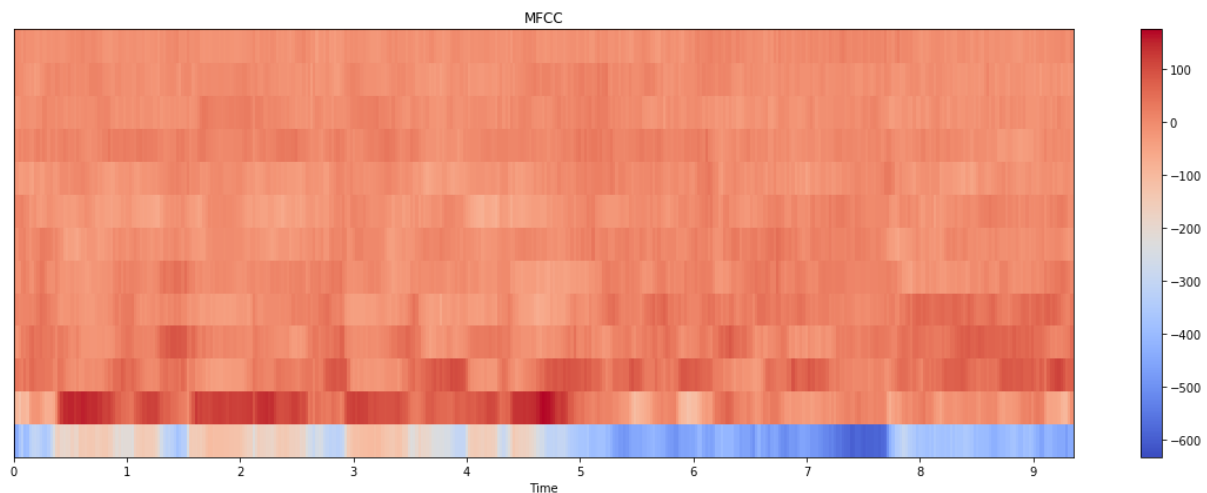
Afin d'obtenir le FFT d'une fenêtre, on doit premièrement créer un spectrom, le spectrom c'est comme une représentation des changements de l'amplitude de signal en fonction de leur fréquence



puis en faire une rotation par 90° et ce dernier (ce n'est pas d'une rotation mais comme l'inverse de la représentation précédente), après cette étape, nous commençons à attribuer pour chaque fréquence, une couleur spécifique, on peut utiliser soit gray scale [0,255] c'est-à-dire 256 possibilités, ou RGB[0,255] à la puissance trois (pour le rouge, le vert, le bleu, c'est-à-dire $256 \times 256 \times 256$, la fréquence élevée est représentée par une couleur foncée. Avec gray scale :



Avec RGB :



Après qu'on crée notre fft pour chaque fenêtre, on combine toutes ces FFT, et on obtient à la fin notre spectrogramme.

Le MFCC vs Mel Spectrogram :

1. mel Spectrogram :

Def : Un Mel spectrogram est une représentation visuelle du spectre des fréquences dans un signal audio au fil du temps. Il affiche l'intensité des différentes composantes de fréquence dans le signal audio, utilisant généralement une palette de couleurs pour représenter l'intensité. Il fournit une représentation en 2D des données audio, où l'axe des x représente le temps, l'axe des y représente la fréquence (généralement sur l'échelle de Mel), et l'intensité de couleur représente l'amplitude ou l'énergie de chaque composante de fréquence à différents intervalles de temps.

- MFCC

Def : Les coefficients cepstraux en fréquence de Mel (MFCC, pour Mel-Frequency Cepstral Coefficients) sont un ensemble de coefficients qui capturent les caractéristiques spectrales d'un signal audio. Ils sont dérivés du spectrogramme de Mel mais sont ensuite traités davantage pour extraire des informations pertinentes. Les MFCC sont généralement utilisés comme caractéristiques pour des tâches telles que l'analyse de la parole et du son. Ils capturent des informations sur la forme de l'enveloppe spectrale, ce qui peut être utile pour des tâches telles que la reconnaissance vocale.

Bien que les spectrogrammes de Mel et les MFCC soient dérivés de la même échelle de Mel, ils servent des objectifs différents. Le spectrogramme de Mel fournit une représentation temps-fréquence du signal audio, tandis que les MFCC sont une représentation compacte des caractéristiques spectrales.

Les chiffres qu'ils génèrent sont différents en raison de ces objectifs distincts et des étapes de traitement.

Image MFCC et Mel spectrogram

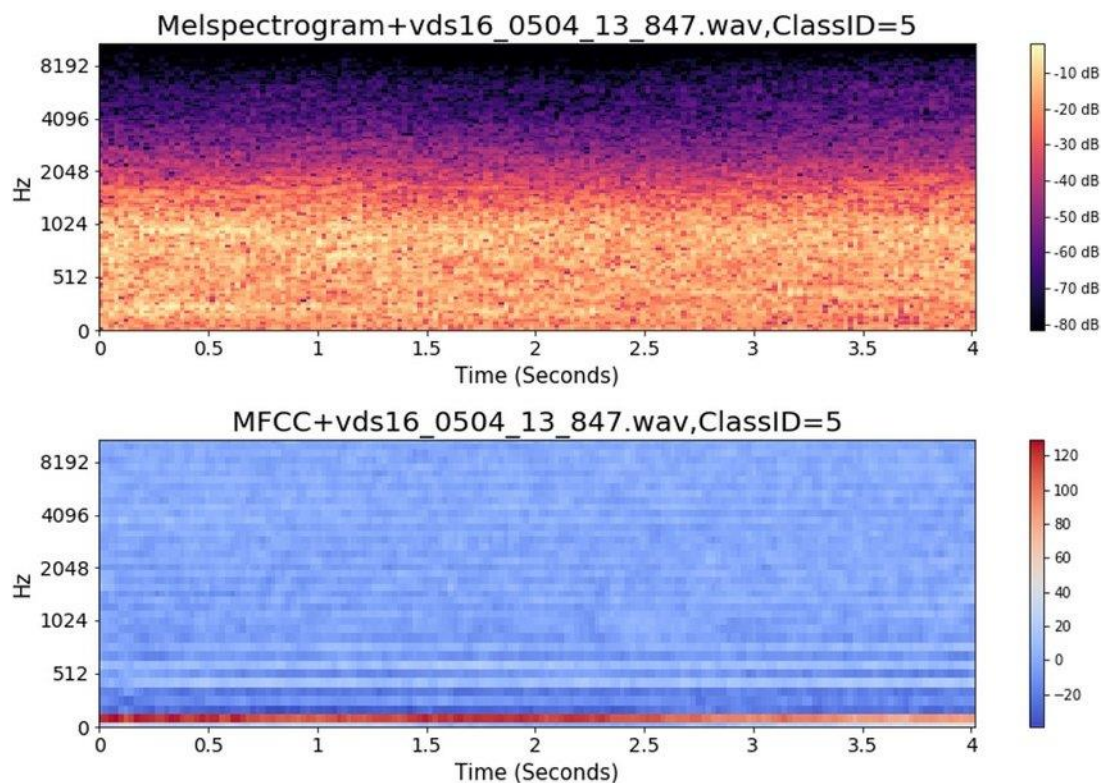
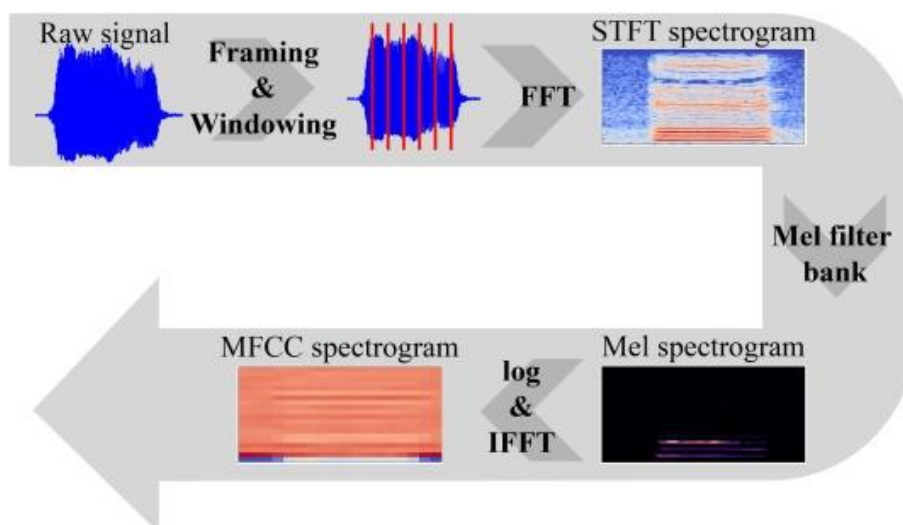


Image : representation de la demarche utilisé pour passer d'un signal audio vers un MFCC :



Ensuite l'application de Neural Network sur ces données préparées.

Partie 2 : traitement de text

Pour la partie textuelle, il faut appliquer uniquement des traitements de NLP (Natural Language Processing) sur les phrases. Ces traitements comprennent principalement la tokenization et la lemmatisation :

Tokenization : c'est le fait de diviser une phrase en unités discrètes appelées tokens. Ces tokens peuvent être des mots individuels, des ponctuations, ou d'autres éléments significatifs tels que des chiffres ou des symboles. La tokenization permet de structurer le texte en éléments plus petits et facilite ainsi son analyse ultérieure.

Lemmatization : c'est le fait de réduire des mots à leur forme canonique ou de base, appelée lemme. Cela implique de transformer les mots en leur forme principale, afin de normaliser le texte et de réduire la dimensionnalité du vocabulaire. Par exemple, les mots "manger", "mangeant" et "mangé" seraient tous réduits à leur lemme "manger".

Après ça, on a besoin d'appliquer un modèle de ML, le modèle que j'ai choisi dans mon cas est le Naïve Bayes classifier,...

Partie 3 : résumé automatique

dans cette partie il n'y a pas de données à préparer, le modèle et l'implémentation de la manière de traitement et le seul besoin pour cette partie, le modèle choisi pour cette détection est TF-IDF, c'est un modèle basé sur le calcul de la fréquence des mots qui composent le document, afin de prendre uniquement les phrases importantes et les combiner pour obtenir le résumé.

Partie 4 : l'interface graphic

il existe deux étapes à faire, la première est de créer une interface utilisateur et implémenter le côté client avec ajax.

La deuxième étape est de créer un serveur node.js dans lequel on définit les routes et utilise child_process pour créer des sous routes où on va faire l'exécution de nos modèles réalisés avec python. (après l'enregistrement des

moel ,en créé un code pythone dans le quelle en faire le loading de notre model et l'executé).

3.2 Mis à jour

Après une longue période de travail, j'ai rencontré des problèmes liés au traitement des fichiers audio. La génération des spectrogrammes MFCC à partir des fichiers audio .WAV a pris 12 heures d'exécution et nécessité 14 Go d'espace de stockage pour les spectrogrammes. De plus, la création du CNN a nécessité l'utilisation de TensorFlow, mais un module appelé distutils manquait. Après des recherches, j'ai découvert que ce module avait été supprimé de la dernière version de Python, la 3.12. J'ai donc essayé de passer à la version 3.11 où le module était encore disponible, mais j'ai rencontré une erreur lors de l'exécution du programme : "ImportError: DLL load failed: Une routine d'initialisation d'une bibliothèque de liens dynamiques (DLL) a échoué." J'ai essayé diverses méthodes pour résoudre ce problème, mais en vain. Finalement, j'ai décidé de poursuivre mon travail sans utiliser les fichiers audio, en me concentrant uniquement sur la détection de la langue et la résumé automatique de texte.

4. Planning et organisation du travail

4.1. Planning prévisionnel

Voici un tableau qui représente les taches, leurs durées et les dépendance entre ces taches :

Etudes sur le NLP et les technologie basic de machine learning T1	1 week	
etudes de base sur le langage de programmation python T2	1 week	
etudes de la bibliotheque NLTK et pandas T3	4 days	T2
etude de la bibliotheque sklearn T4	1 week	T2
Choisir l'algorithme de ML à appliquer pour la detection de langue. T5	1 week	T1
recherche de la base de données convenable pour la detection de langue T6	1 day	T5
Création du model de detection de langue T7	3 days	T2-T3-T4-T5-T6
Choisir l'algorithme de ML pour le résumé automatique T8	1 week	T8-4-3-2
Création de model de résumé automatique T9	3 days	
Création de l'interface graphique avec html et css(bootstrap) T10	2 hours	T10
Réalisation de côté client avec ajax T11	5 hours	
Réalisation des routes au cours de l'implémentation de serveur node.js T12	2 hours	T7,T9
Création des scripts python pour exécuter les modèles enregistrés T13	1 hour	T12,T7
La liaison entre le serveur et le model de detection de langue T14	2 days	T12,T9
La liaison entre le serveur et le model de résumé automatique T15	3 days	T7,T9
Teste des modèles T16	1 hour	T12,T13,T14,T15
Test de la connection entre le serveur et les modèles T17	1 day	T7,T9,T10,T11,T12,T13,T14,T15

Échéances :

La date au plus tard de chaque tâche représente la date de départ de la dernière tâche qui dépend de cette dernière, plus d'information dans le diagramme de GANTT dans la partie conception

Suivi et mise à jour :

Plus d'information dans le diagramme de GANTT dans la partie conception

5. livrables attendus

Le résultat attendu, contenir un ensemble des composants regroupés pour fournir par la suite une interface qui sert à détecter la langue d'un texte fournie et faire un résumé automatique.

Concernant le modèle de détection de langue, il est implémenté avec l'algorithme Naïve Bayes Classifier .

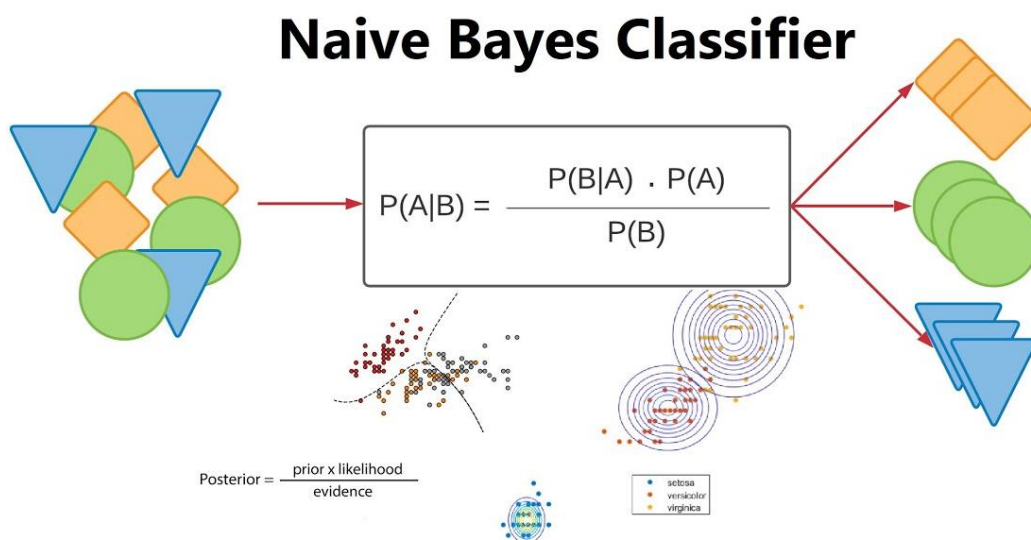
C'est quoi le modèle Naive Bayes Classifier :

D'après wikipedia :

``La classification naïve bayésienne est un type de classification bayésienne probabiliste simple basée sur le théorème de Bayes avec une forte indépendance (dite naïve) des hypothèses. Elle met en œuvre un classifieur bayésien naïf, ou classifieur naïf de Bayes, appartenant à la famille des classifieurs linéaires.``

Naive Bayes est un algorithme de classification simple et rapide, parfaitement adapté à des ensembles de données volumineux. Il est largement utilisé dans divers domaines tels que le filtrage de spam, la classification de texte, l'analyse de sentiment et les systèmes de recommandation. Son fonctionnement repose sur le théorème de Bayes, qui permet de calculer la probabilité d'appartenance à une classe inconnue à partir des probabilités conditionnelles des caractéristiques observées. Sa simplicité et son efficacité en font un choix populaire pour de nombreuses applications de traitement de données.

Voici une image explicative de Naive Bayes Classifier



Voici un exemple explicatif :

L'algorithme Naive Bayes est utilisé pour prédire si un e-mail est du spam ou non-spam. Il utilise une méthode simple basée sur la probabilité pour cette prédiction.

Imaginez qu'on a une boîte de réception remplie d'e-mails. Certains de ces e-mails sont des spams, tandis que d'autres sont des e-mails légitimes.

Pour entraîner le modèle Naive Bayes, nous examinons un grand nombre d'e-mails déjà étiquetés comme spam ou non-spam. Nous regardons les mots qui apparaissent dans ces e-mails et calculons la probabilité qu'un e-mail soit du spam en fonction de ces mots.

Par exemple, si nous remarquons que les mots "gratuit", "offre", et "promo" apparaissent souvent dans les e-mails étiquetés comme spam, nous attribuons une probabilité plus élevée à un e-mail contenant ces mots d'être du spam. Ensuite, lorsque nous recevons un nouvel e-mail, le modèle Naive Bayes examine les mots qu'il contient et calcule la probabilité que cet e-mail soit du spam ou non. Il prend en compte toutes les probabilités des mots et combine ces informations pour obtenir une probabilité globale.

Si la probabilité que l'e-mail soit du spam dépasse un certain seuil, alors le modèle le classe comme spam. Sinon, il le classe comme non-spam.

L'avantage de ce modèle c'est qu'il s'implémente et à comprendre, aussi il ne nécessite pas une large quantité de données pour un résultat exact et cela car en raison de son approche probabiliste simple et efficace, il suppose que les caractéristiques sont indépendantes les unes des autres, ce qui signifie qu'il ne nécessite pas autant de données pour estimer les probabilités conditionnelles des différentes caractéristiques.

J'ai essayé d'utiliser le modèle logistique, mais même si l'accuracy est élevée 0.95, il me donne un résultat faussé, voici une capture de démonstration :

```
from sklearn.linear_model import LogisticRegression
ada = LogisticRegression()
y_pred = ada.fit(X_train,y_train).predict(X_test)
```

✓ 28m 34.8s

```
from sklearn.metrics import f1_score, confusion_matrix
f1 = f1_score(y_test, y_pred, average='macro')
```

✓ 0.2s

```
print(f'Accuracy = {ac:.2f}')
```

✓ 0.0s

Accuracy = 0.95

```
def prediction(text):
    x = cv.transform([text]).toarray()
    lang = ada.predict(x)
    lang = le.inverse_transform(lang)
    print('The language is in', lang[0])
```

✓ 0.0s

```
prediction('Bonjour maman')
```

✓ 0.0s

```
prediction('Bonjour maman')
```

✓ 0.0s

• The language is in Russian

La deuxième étape et pour le modèle de résumé automatique, le modèle que j'ai choisi est le Tf-IDF

C'est quoi TF-IDF?

Abréviation de (Term Frequency-Inverse Document Frequency) C'est un modèle de Machine Learning basé sur une technique de NLP, utilisé pour évaluer l'importance des différents mots dans une phrase.

Le modèle TF-IDF est utilisé pour évaluer l'importance des mots ou des phrases dans un document afin de sélectionner les informations les plus pertinentes à inclure dans le résumé final. Il existe plusieurs mesures pour déterminer la fréquence de mot (TF), soit le calcul du nombre d'apparition de mot, ou faire une division par rapport au nombre des mots dans le document,

par contre IDF représente l'inverse de la fréquence des documents (Inverse Document Frequency). Cette mesure évalue l'importance d'un terme dans l'ensemble du corpus de documents. Les termes qui apparaissent dans de nombreux documents ont une faible valeur IDF, mais les termes qui apparaissent dans peu de documents ont une valeur IDF élevée, et sont considérés comme les plus informatifs.

A l'aide de ça, il évalue l'importance des mots ou des phrases dans le document et sélectionne les informations les plus pertinentes à inclure dans le résumé final en multipliant le TF et le IDF. Les termes avec les poids TF-IDF les plus élevés sont généralement sélectionnés pour représenter les points clés du document d'origine dans le résumé.

J'ai choisi ça puisque c'est un modèle car il est le plus adapté pour le résumé automatique grâce à plusieurs raisons, la première c'est qu'il est adapté pour la variété des domaines des documents, aussi il ne s'intéresse pas uniquement aux mots les plus fréquents dans le texte mais aussi aux rares, de plus il est facile à implémenter, ne nécessite pas de modélisation complexe ou d'apprentissage supervisé, ce qui peut rendre le processus plus rapide et plus efficace.



Le TF représente la fréquence de mot par rapport au document de texte, et le IDF représente la fréquence de mot par rapport à chaque partie de document (texte).

$$TF(t, d) = \frac{\text{number of times } t \text{ appears in } d}{\text{total number of terms in } d}$$

$$IDF(t) = \log \frac{N}{1 + df}$$

$$TF - IDF(t, d) = TF(t, d) * IDF(t)$$

L'aspect mathématique de (TF-IDF).

Maintenant, l'interface utilisateur :

Pour l'interface graphique, il doit contenir uniquement une zone de texte l'une avec l'écriture désactivée, puis un bouton, l'une pour l'action de détection de langue et l'autre pour la summarisation.

Sur le bouton, il y a un listener sur lequel le code ajax est créé pour envoyer les données nécessaires au serveur et retourner le résultat d'après ce dernier.

Puis, pour le côté serveur, les routes sont définies, notamment deux routes demandées, les deux pour la méthode post, la première pour le bouton get language, il doit retourner les données envoyées après la click sur le bouton (code implémenté avec ajax), puis interagit avec le modèle enregistré de détection de la langue, et envoie le résultat du modèle afin de l'afficher pour l'utilisateur.

La deuxième et pour le bouton get summary, dans cette dernière il faut prendre en considération si le bouton est getLanguage et cliquer ou non, car la fonction qui représente le modèle nécessite deux paramètres le premier et la langue du texte inséré et le deuxième est le texte, si le bouton est déjà cliqué alors la langue obtenue de la zone où la langue est affichée, sinon, il faut faire l'appel au modèle de détection de langue, puis obtenir la langue et l'utiliser comme paramètre avec le texte pour obtenir le résumé automatique.

Et enfin, le test final pour valider le résultat final.

Conclusion

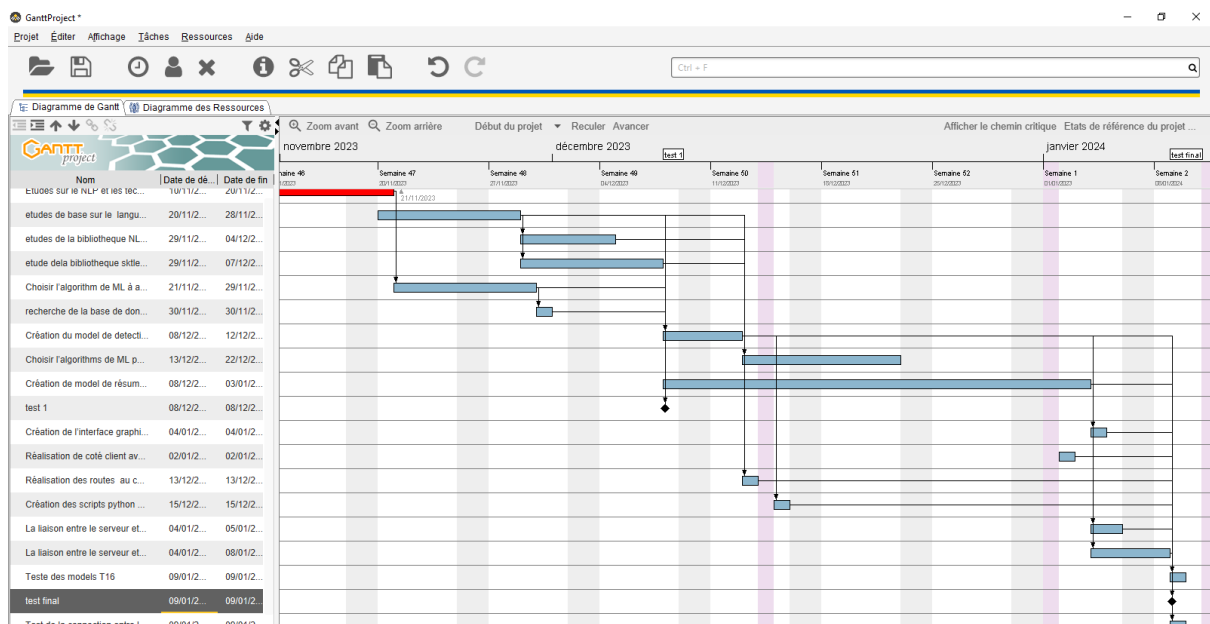
je souligne ,comme conclusion que le projet vise à répondre à un besoin essentiel dans le domaine du traitement automatique du langage naturel en fournissant une solution efficace pour la détection de la langue et le résumé automatique de texte. En utilisant des technologies telles que le Naive Bayes classifier et le modèle TF-IDF

j'ai essayé de bien planifier le projet ,et determiner les taches principal afin de ne pas raté la durée de chacune des taches depuis l'exploration des technologies de base jusqu'à la mise en œuvre des modèles et la création de l'interface utilisateur.

Je reste déterminé à poursuivre mes efforts et à livrer un rapport complet et satisfaisant, répondant aux attentes de toutes les parties prenantes impliquées.

II. Conception du projet

Image représente le diagramme de gantt de projet :



Imag represente le diagramme de pert de projet :

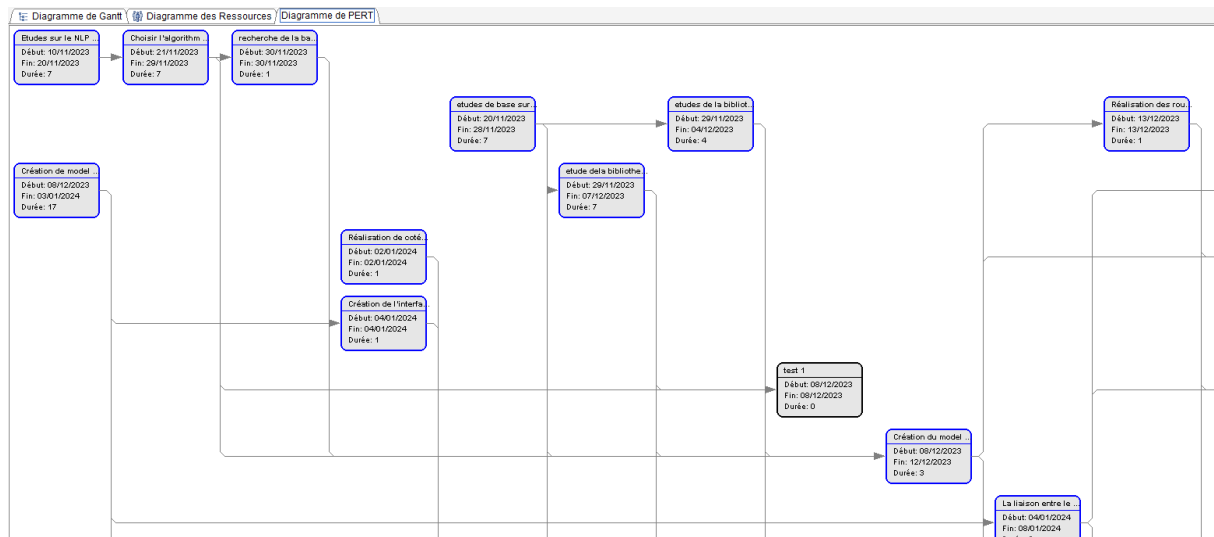
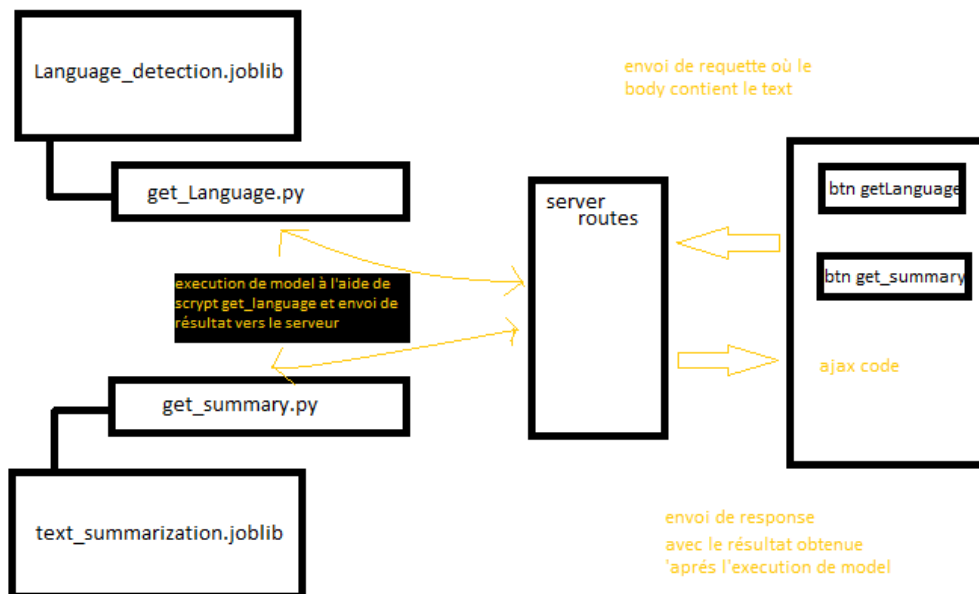


Image représente une simulation des different interaction des composant de plateforme :



III. Outils est technologies utilisé

La réalisation de ce projet naicessite une large gamme des outils et technologie :

Environnement de développement intégré (IDE) :

J'ai pas essayé d'utiliser un différent IDE que visual studio code puisque c'est le plus populaire, la seule chose que j'ai ajoutée sont les extensions, la première est de python, suivie par l'extension de jupyter afin d'obtenir un environnement similaire de jupyter notebook pour les fichiers avec l'extension .ipynb

Un Jupyter Notebook est un environnement de développement interactif qui permet d'écrire et d'exécuter du code en Python (ou dans d'autres langages comme R ou Julia) de manière progressive, combinant des cellules de code, de texte, d'images et de visualisations dans un seul document.

Gestionnaire de packages :

Pour le gestionnaire des packages, j'ai utilisé pip puisqu'il est le gestionnaire de packages standard pour Python.

Bibliothèques Python :

Il existe nombreuses bibliothèques utilisées dans mon projet, j'ai déjà les mentionnées dans la partie de 1.4 de Chapitre I 'Etat actuel des ressources et fonctionnement du préexistant, mais je les rementionne avec une brève explication :

NumPy : utilisé pour le calcul numérique.

pandas : utilisé pour la manipulation et l'analyse des données.

Matplotlib: utilisé pour la visualisation des données.

scikit-learn : utilisé pour l'apprentissage automatique.

NLTK (Natural Language Toolkit) ou spaCy : utilisé pour le traitement du langage naturel (NLP).

Environnement virtuel :

L'environnement virtuel est l'outil principal qui permet d'installer les packages Python spécifiques pour un projet sans interférer avec les autres projets ou le système global. Dans mon cas j'ai utilisé conda (Anaconda), conda est un outil polyvalent, il est pas utilisé uniquement pour l'environnement virtuel mais aussi pour la gestion des packages et plus, j'ai choisie de l'utiliser car il gère

automatiquement les dépendances entre les packages, ce qui garantit que toutes les dépendances requises sont correctement installées et compatibles, aussi c'est facile à utiliser et simplifie la création, la gestion et l'utilisation des environnements virtuels avec quelques commandes simples.

Les algorithmes de ML et les techniques de NLP :

NaiveBayesClassifier : Le Naive Bayes Classifier est un algorithme de classification supervisée simple mais puissant, souvent utilisé pour la classification de texte, telle que la détection de spam, l'analyse de sentiments, ou la catégorisation de documents. Il est basé sur le théorème de Bayes et suppose que les caractéristiques sont indépendantes les unes des autres.

TF-IDF : TF-IDF est une méthode utilisée pour évaluer l'importance d'un terme dans un document par rapport à un corpus de documents. Elle prend en compte à la fois la fréquence d'apparition d'un terme dans un document (TF) et l'inverse de sa fréquence d'apparition dans l'ensemble des documents du corpus (IDF).

Vous pouvez voir plus d'information dans la partie 4 chapitre 1 (livrables attendue).

HTML : c'est le langage principal pour le côté statique des sites web.

Bootstrap : c'est un framework open source développé qui fournit une collection d'outils et de composants front-end prêts à l'emploi pour le développement rapide et la conception de sites web et d'applications web réactives, il permet de personnaliser facilement l'apparence de site avec les thèmes prédéfinis.

Ajax (Asynchronous JavaScript and XML) : c'est une technique de développement web qui permet de créer des applications web interactives et dynamiques en utilisant JavaScript et des requêtes HTTP asynchrones, c-à-d mises à jour asynchrones et partielles des contenus web sans rechargement complet de la page.

Node.js : un environnement d'exécution JavaScript côté serveur, construit sur le moteur JavaScript V8 de Chrome, permet d'exécuter du code JavaScript côté serveur, et utiliser pour créer des applications serveur, des API, des applications web en temps réel .. .

child_process : un module intégré à Node.js qui permet d'exécuter des processus externes depuis une application Node.js. et fournit plusieurs

méthodes pour créer, contrôler et communiquer avec ces processus externes, dans notre cas on s'intéresse par la méthode `spawn` : qui est utilisée pour démarrer un processus externe et obtenir un flux de sortie (`stdout`) et un flux d'erreur (`stderr`) pour l'exécution de notre script de l'exécution des modèles python créés.

Express : Express est un framework web minimaliste et flexible pour Node.js, qui simplifie le processus de création d'applications web et d'API. Il fournit un ensemble de fonctionnalités robustes pour la gestion des routes, la gestion des requêtes et des réponses http.

Body-parser : un middleware pour Express qui permet de récupérer les données envoyées dans le corps des requêtes HTTP, telles que les données de formulaire ou les données JSON. Il analyse le corps des requêtes et le transforme en un objet JavaScript utilisable par l'application.

IV. Réalisation

l'environnement de travail :

Premièrement, j'ai commencé par l'installation de python, j'ai travaillé avec la version 3.11 puisqu'il contient des modules pas encore supprimés par rapport au v3.12, puis j'ai téléchargé anaconda afin d'avoir un environnement virtuel de python, puis j'ai ajouté pip et python dans Path avec les variables d'environnement car j'ai rencontré un problème de commande introuvable pour pip et python, puis j'ai téléchargé l'extension python et jupyter pour obtenir un environnement similaire que jupyter notebook.

modèle de détection de langue :

Au début, j'ai commencé par l'ajout d'un fichier csv qui contient les données que je serai utilisé pour la partie de training et pour la partie de testing, puis j'ai commencé à importer les packages nécessaires pour la réalisation de ce modèle, qui sont pandas, numpy, nltk, re

Puis j'ai chargé les données à partir du fichier csv et j'ai passé vers la préparation de cette dernière, la première étape et de supprimer les caractères spéciaux et les chiffres, et en mettre tout le texte en minuscule, après j'ai encodé

les langue en valeurs numériques,et j'ai utilisé ces valeurs comme étiquette , chaque ligne dans le tableau qui contient les données sera associée à la valeur numérique de la langue de sa colonne correspondante comme étiquette,et comme ça ,j'ai créé plusieurs classes,chacune pour un langage spécifique.

L'étape suivante est de convertir les phrases en vecteurs,et cela basé sur le CountVectorizer(),premierement CountVectorizer commence par construire un vocabulaire "Bag-of-Words,, à partir de toutes les mots uniques dans l'ensemble des phrases. Chaque mot unique est considéré comme une "caractéristique" du texte,puis il attribue un index à chaque mot du vocabulaire ,et comme cela,chaque phrase est convertie en vecteur dont les nombres sont les index attribuer au mots qui composent cette phrase,la même procédure est appliquée à toutes les phrases de la base de données.

Après cette étape ,les données sont préparées ,je les divise en données de training et données de testing,pour la partie training 80% des données sont fournies ,par contre pour la partie testing ,uniquement 20% sont fournies,j'ai utilisé cette démarche car Plus l'ensemble d'entraînement est grand, plus le modèle aura de données sur lesquelles apprendre, ce qui peut donner comme résultat une meilleure performance du modèle .

Après j'ai initialisé la classification Multinomial de Naive BayesClassifier ,et j'ai entraîné ce modèle avec les vecteurs des phrases de training et les indices des langues avec la méthode fit, je passe vers l'évaluation sur l'ensemble de test

Sur les 20% des données que j'ai laissées pour la partie de testing ,et pour assurer que chaque classe a le même poids dans le calcul du score F1 global,j'ai utilisé `average='macro'`,et comme cela les classes plus petites seront avoir le même poids que les classes plus grandes dans le calcul du score global (il sera moyenné pour obtenir un score global) , et j'affiche à la fin le score ,puis j'enregistre le modèle sous à l'aide de dump.

Test de modèle de détection de langue :

Concernant la partie testing,il est inclus dans la réalisation de projet (l'étape avant dernière),mais après j'ai essayé de tester avec des différentes phrases ou les mots sans les mêmes pour les deux langues,ce que j'ai remarqué que si je donne un mot qui est le même pour deux langues (anglais et autre langue) sans article,ce dernier sera considéré comme anglais ,le cas que j'ai donné n'a pas une réponse exacte mais aussi la réponse donnée par le modèle reste correct puisqu'il n'a pas d'article pour spécifier.

model de résumé automatique

Pour ce model,j'ai choisie d'utilisé la methode TF-IDF,j'ai déjà parlé de cette dernière dans la partie 4 chapitre 1 (livrable attendus).

Premièrement je faire le stemming pour réduire les mot a leur format de base et éliminé les stop words ,aussi ignore les mot avec un nombre de caractère inférieur à 3 puisqu'il ne serons pas contient des info de base que je cherche dans le documment ou le text ,puis je créé une table de fréquence de ce mot en fonction des phrases qui compose le text,j' utilise la methode de NLTK sent_tokenize pour capté les phrases de text et je réalise une matrice de fréquence des mot pour chaque phrase,ensuite je pass vers le calcule de TF,avec la matrice des mot pour chaque phrase que j'ai fait je calcule la matrice TF (Term Frequency) en divisant chaque fréquence par le nombre total de mots dans la phrase,je garde ce calcule et je pass vers le calcule de nombre de phrases contenant chaque mot ,et cela a l'aide de la matrice de fréquece des mot pour chaque phrase que j'ai déjà réalisé,puis je pass vers le calcule de l'IDF,l'iDF ,et par la suite je calcule le score en fesosns une somme de TF et IDF , et j'utilise ce résultat pour calculer le score des phrases ,ensuite je calcule le scors moyenn des phrases et je determine un seuil, le seuil represente le score moyen des phrases multiplier par un facteur spécifique ,dans mon cas j'ai utilisé 1.15 puisqu'il est recommandé,a la fin ,pour le résumé automatique, en choisissant uniquement les mot avec un score supérieur à ce seuil ,en combinant ces phrases entre eux et en obtenir par l suite un résumé de notre texte /documment .

test de model do résumé automatique

J'ai testé avec deux fichiers textuelles et j'ai obtenue un résultat satisfant ,dans laquelle le nombres des mot et réduit ,et le text donne une idée général de l'ancienne documment détaillé.

création de l'interface graphic

Pour l'interface graphic,j'ai utilisé premierement htmlml pour créer deux zone de text la premier j'ai gardé la notion de l'écriture activé mais je la supprimé de la de l'autre, aussi deux bouton une pour le résumé et l'autre pour la langue ,et une label 'langue, suivie par une balize p vide ou je serai inséré la langue détecté.

Ensuite j'ai utilisé seulement bootstrap pour donner une taille grande au zone de text et un style pour les bouton ,et alignement au centre.

implémentation de coté Serveur(Node.js)

Au coté serveur,j'ai réalisé deux routes ,une pour le résumé et l'autre pour la detection de la langue.

Premièrement j'ai importé tous ce qui est naicessaire,les moduls,les middlewear,et ainsi de suite,puis j'ai pasé vers la création des routes .

- Pour la route detection de langue('/getLanguage') :

C'est une route post avec le chemin déjà mentionné,premierement j'ai Récupéré le texte envoyé dans le corps de la requête POST à l'aide de req.body.text, puis j'ai Créé un processus enfant dans laquelle j'exécute le script Python language_detection.py,. ce script import le model que j'ai déjà enregistré dans un fichier .joblib et exécute les methodes correspond a ce model (j'ai fais la même chose pour le résumé). L'execution de ces script est fait à l'aide de la méthode spawn du module child_process que j'ai importé précédament et j'envoie le texte récupéré dans le corps de la requête via stdin pour qu'il puisse utilisé et je termine l'écriture dans l'entrée standard du processus Python. Et je commence à écouter les données envoyées par le processus Python sur sa sortie standard qui représentent dans ce cas la langue de texte fournie,

Et en fin Envoie une réponse JSON au client(implémenté par ajax) contenant la langue détectée.

- Pour la route de résumé ('getSummary') :
- Aussi c'est une route post avec le chemin que j'ai mentionné, premièrement je fais la feth ala route getLanguage afin d'obtenir la langue de ce text fournie puis je crée un processus enfant pour exécuter le model de summarisation et j'envoie le texte récupéré de body de requête plus la langue détecté après l'appel de la route getLanguage,puis j'obtenir mon résultat de l'exécution de model et j'envoie la réponse au client qui est dans ce cas un summary de text.

test de coté serveur

J'ai utilisé une extension dans vs code qui me permet de généré des requêtes ,j'ai défini le port de l'écoute et l'adress localhost et j'ai essayé avec un text,et heureusement ,j'ai obtenue un response qui contient dans le body la langue détecté.

J'ai fait la même chose pour la partie de résumé automatique.

implémentation e coté client (Ajax)

Pour cet partie ,j'ai déjà donné des id au deux boutons et au zones de texte ainsi que le div ou la langue détecté sera affiché ,je créé des variable pour obtenir ces composant d'après ces id ,puis j'utilise la methode `addeventListener` de l'évènement `click` ,et j'ai défini la fonction à exécuté ,pour les deux boutons (`get language` et `get Summary`),la seul chose que je serai fait et de préparé une requette `html` vers le serveur `node.js` où j'ai déjà définie les deux routes,cette requette utilise la méthode `post` et envoi dans le `body` le texte fournie par l'utilisateur dans la zone de texte.

Voici une capture d'écran pour visualisé comment ça marche la détection et la génération de résumé :

L'interface :



The screenshot displays a web interface with two main sections: 'Input Text' and 'Summary'. The 'Input Text' section contains a large text area with the placeholder text 'Enter your text here...'. Below this text area are two blue buttons: 'Get Language' and 'Get Summary'. Below the buttons, the text 'Language:' is visible. The 'Summary' section is an empty rectangular box on the right side of the interface.

Test avec un texte en allemand :

Détection de langue :

Input Text

Die deutsche Geschichte ist reich an Ereignissen und bedeutenden Persönlichkeiten. Sie erstreckt sich über Jahrhunderte und umfasst eine Vielzahl von Epochen und Entwicklungen. Ein wichtiger Meilenstein war die Gründung des Heiligen Römischen Reiches im Jahr 800 n. Chr. durch Karl den Großen, der als einer der bedeutendsten Herrscher des Mittelalters gilt.

Get Language

Get Summary

Language:

German

Summary

Génération de résumé :

Input Text

Die deutsche Geschichte ist reich an Ereignissen und bedeutenden Persönlichkeiten. Sie erstreckt sich über Jahrhunderte und umfasst eine Vielzahl von Epochen und Entwicklungen. Ein wichtiger Meilenstein war die Gründung des Heiligen Römischen Reiches im Jahr 800 n. Chr. durch Karl den Großen, der als einer der bedeutendsten Herrscher des Mittelalters gilt.

Im Laufe der Jahrhunderte erlebte Deutschland eine Vielzahl von

Get Language

Get Summary

Language:

German

Summary

In der Musik haben Komponisten wie Ludwig van Beethoven, Johann Sebastian Bach und Wolfgang Amadeus Mozart unvergessliche Meisterwerke geschaffen, die bis heute als Höhepunkte der klassischen Musik gelten.

V. Conclusion

Dans le cadre de mon projet de fin d'études, j'ai développé une plateforme de détection de langue et de résumé automatique de texte dans le but de faciliter la compréhension et l'analyse de documents dans différentes langues. Mon objectif principal était de créer un outil polyvalent capable de traiter les textes provenant de sources diverses et de générer des résumés concis et informatifs.

Pour atteindre cet objectif, j'ai utilisé une approche basée sur le traitement automatique du langage naturel (NLP) et l'apprentissage automatique (ML). J'ai commencé par collecter des données textuelles dans différentes langues, puis j'ai procédé à un prétraitement rigoureux, comprenant le nettoyage des données, la tokenisation et la suppression des stopwords.

En ce qui concerne la détection de langue, j'ai développé un modèle basé sur des caractéristiques linguistiques telles que la fréquence des mots et la structure syntaxique. Pour le résumé automatique, j'ai adopté une approche basée sur la fréquence des termes-inverse du document (TF-IDF) pour identifier les phrases les plus importantes dans le texte. En utilisant cette méthode, j'ai pu générer des résumés pertinents en préservant l'essence du contenu original.

Bien que mon système ait montré des performances encourageantes, il présente encore certaines limitations. Par exemple, il peut rencontrer des difficultés avec les langues moins représentées dans mes données d'entraînement, et la qualité des résumés peut varier en fonction de la complexité du texte source. De plus, la détection de langue ne couvre pas toutes les langues du monde, mais seulement quelques-unes selon la base de données disponible.

De plus, je n'ai pas pu aller plus loin dans le traitement des audio en raison d'un problème lié au package TensorFlow qui m'a pris beaucoup de temps, mais malheureusement sans résultat.

Pour les futures directions de recherche, j'envisage d'explorer des techniques de ML pour améliorer la précision de la détection de langue, ainsi que d'intégrer la notion de traitement des audio en utilisant des méthodes de deep learning et de réseaux neuronaux. De plus, je souhaite améliorer la qualité des résumés automatiques et découvrir d'autres techniques à l'avenir.

Ce projet a représenté une étape importante dans ma carrière dans le domaine du traitement automatique du langage naturel. Je tiens à exprimer ma gratitude à mon encadrant pour son soutien et ses conseils tout au long de ce projet, ainsi qu'à l'ESTF pour m'avoir donné cette opportunité d'apprentissage et de développement.

Annexes :

https://fr.wikipedia.org/wiki/Classification_na%C3%AFve_bay%C3%A9sienne

<https://www.oracle.com/ch-fr/artificial-intelligence/machine-learning/what-is-machine-learning/>

https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html

[https://www.learndatasci.com/glossary/tf-idf-term-frequency-inverse-document-frequency/#:~:text=Using%20scikit%2Dlearn-What%20is%20TF%2DIDF%3F,%2C%20relative%20to%20a%20corpus\).](https://www.learndatasci.com/glossary/tf-idf-term-frequency-inverse-document-frequency/#:~:text=Using%20scikit%2Dlearn-What%20is%20TF%2DIDF%3F,%2C%20relative%20to%20a%20corpus).)

<https://www.ibm.com/topics/natural-language-processing>

<https://www.analyticsvidhya.com/blog/2021/04/role-of-machine-learning-in-natural-language-processing/>

<https://ieeexplore.ieee.org/document/9955539>

<https://fr.mathworks.com/matlabcentral/answers/597799-difference-between-fft-and-pspectrum>

<https://www.science.gov/topicpages/f/fft+spectrum+analyzer.html>