

EFREIMOTION

MASTER CAMP

ANAS BOUJATOUY, MORDJANE ARIBI, FELIX DEVYNCK, FODE DIEDHIOU
GABRIEL DE CHAMBOST, RAYEN HADDAD

DATA SCIENCE | GROUPE 225



PARIS PANTHÉON-ASSAS UNIVERSITÉ

Sommaire

Introduction.....	2
Contexte et motivation.....	2
Objectifs du projet.....	2
Méthodologie.....	3
Organisation	3
Collecte et préparation des données	3
Choix de l'algorithme de machine Learning.....	4
Entraînement du modèle	5
Résultats et Discussion	5
Développement.....	5
Application.....	6
Performance de l'algorithme.....	7
Comparaison avec d'autres approches existantes	8
Limitations et améliorations possibles	9
Conclusion	9
Contributions du projet.....	9
Réflexion sur les résultats obtenus.....	10
Perspectives.....	10

Introduction

Contexte et motivation

Lors de notre troisième année à EFREI Paris, dans le cadre du Mastercamp, nous avons été confrontés à un projet d'envergure lié à notre future filière : DATA SCIENCE. Nous avons eu le choix parmi une sélection de trois projets prédéfinis par nos enseignants. Après avoir étudié attentivement chacun d'entre eux, nous avons choisi de travailler sur le premier projet proposé, à savoir l'analyse fine d'avis d'utilisateurs.

L'analyseur de sentiments est une application ou un algorithme capable de comprendre et d'interpréter les émotions et les opinions exprimées dans des textes ou des messages. Il s'agit d'un domaine d'étude omniprésent et en constante évolution, car il permet d'extraire des informations précieuses à partir de données textuelles. Nous pouvons le retrouver dans de nombreux cas, notamment sur internet et sur les réseaux sociaux.

Ce projet nous a semblé particulièrement pertinent car il s'inscrit dans un domaine qui touche, de nos jours, toutes les grandes entreprises. En effet, de nombreuses organisations cherchent à comprendre les opinions et les émotions des consommateurs afin d'adapter leurs stratégies de marketing, d'améliorer leurs produits et services, ou encore de gérer leur réputation en ligne. Cette compréhension doit donc permettre aux prestataires de mettre au point une amélioration du produit ou service vendu. L'analyse des sentiments leur permet de tirer des informations précieuses à partir des données textuelles disponibles sur les réseaux sociaux, les forums, les critiques de produits, etc.

La possibilité de développer un outil ou un algorithme capable d'extraire ses informations précieuses à partir de grandes quantités de données textuelles peut être extrêmement gratifiante. Contribuer à la création d'une solution qui a un impact réel sur les décisions commerciales ou la compréhension des opinions des utilisateurs a été une source de motivation pour l'ensemble du groupe.

D'autre part, nous avons vu en ce projet, l'opportunité de mettre en application et de développer nos compétences techniques en programmation et en Machine Learning mais aussi nos compétences humaines en travaillant en équipe et en développant nos compétences commerciales.

Objectifs du projet

Notre principal objectif est de concevoir et de tester un algorithme capable d'analyser avec précision les sentiments exprimés dans des textes. Nous souhaitons exploiter les techniques avancées de Machine Learning, telles que les réseaux de neurones, les Forêts aléatoires ainsi que du one hot encoding pour obtenir des résultats fiables et cohérents.

Pour former notre modèle d'analyseur de sentiments, nous devons rassembler un ensemble de données représentatif et diversifié. Nous nous fixons comme objectif de collecter des échantillons de texte provenant d'une base de données complète avec des produits variés. Pour ce faire, nous avons décidé d'étudier et de créer nos modèles de Machine Learning à partir d'une base de données Amazon.

Nous souhaitons que notre programme soit en mesure à partir d'un texte ou d'une base de données, d'extraire le pourcentage de 5 émotions différentes et d'affiche un nuage de mot pour chaque sentiment afin de permettre aux utilisateurs de comprendre les émotions exprimées dans les commentaires ainsi que les mots liés à ces dernières.

Méthodologie

Organisation

Afin de garantir le bon déroulement de notre projet, nous avons mis en place une organisation précise et efficace en utilisant des méthodologies de gestion de projet. L'approche agile a été au cœur de notre méthodologie, nous permettant de travailler de manière itérative et collaborative.

Nous avons commencé par définir des objectifs clairs et spécifiques, que nous avons ensuite décomposés en tâches plus petites et réalisables. Ces tâches ont été organisées en itérations, avec des délais courts pour favoriser la flexibilité et l'adaptabilité. À chaque fin d'itération, nous avons procédé à des réunions de suivi pour évaluer les progrès réalisés, ajuster les priorités et planifier les étapes suivantes.

L'utilisation d'outils de gestion de projet tels que Trello nous a permis de visualiser et de suivre facilement l'avancement des tâches, d'assigner des responsabilités aux membres de l'équipe et de favoriser la collaboration. De plus, l'utilisation de GitHub nous a permis de gérer efficacement les versions et les modifications du code source, facilitant ainsi le travail d'équipe et le suivi des contributions individuelles.

En adoptant cette approche agile, nous avons pu nous adapter rapidement aux changements et aux imprévus, tout en maintenant une communication fluide au sein de l'équipe. Cela nous a permis d'optimiser notre productivité, de résoudre les problèmes rapidement et d'atteindre nos objectifs dans les délais impartis.

Collecte et préparation des données

Concernant la base de données, nous avons une base de données Amazon disponible en open source sur le site KAGGLE.COM. Notre choix a été motivé par la différence des données qu'elle comporte étant donné la diversité des produits commercialisés par le groupe mais aussi par le fait qu'elle est très complète. En effet, notre base de données possédait, avant tri de notre part, plus de 560 000 lignes et 10 colonnes.

Lien de notre base de données : <https://www.kaggle.com/datasets/arhamrumi/amazon-product-reviews>

Il a ensuite fallu faire la préparation des données. Pour ce faire, nous avons établi les données importantes pour notre projet et réduit la quantité des données pour rendre le lancement de nos programmes plus faciles et rapides.

Nous avons décidé de grader les colonnes suivantes :

- Score : Note sur 5
- Summary : Résumé de l'avis
- Text : Commentaire

Nous avons donc réduit notre nombre de données à 50 000 lignes, nombre suffisant pour pouvoir réaliser des modèles de Machine Learning efficace.

Afin de réaliser nos modèles, nous avons définis la TARGET, c'est à dire la variable cible que l'algorithme doit être capable de prédire. Nous avons donc décidé de créer de nouvelles colonnes possédant un pourcentage pour les 5 sentiments que nous voulons mettre en avant : la joie, la déception, la tristesse, la colère et la surprise.

Afin de pouvoir créer ces colonnes de variable cible, nous avons créé des listes de mots pour chaque sentiment afin de pouvoir réaliser du one hot encoding sur les commentaires. Chacune des listes comporte une centaine de mots en anglais pour s'adapter à notre base de données.

Afin d'adapter le pourcentage à l'ensemble des sentiments, nous avons créé une colonne `sum_pourcentage`.

Aperçu de notre base de données :

Score	Time	Summary	Text	Encoding_joy	Encoding_dis	Encoding_sad	Encoding_ang	Encoding_sur	Joy_pourcent	Disappointme	Sadness_pour	Anger_pource	Surprise_pour	sum_pourcentage
5.0	1303862400.0	good quality c	i have bought	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	0.7194244604	0.0	0.0	0.0	0.7194244604316548	
1.0	1346976000.0	not as adverti	product arrive	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	0.7633587786	0.0	0.0	0.0	0.7633587786259541	
5.0	1267920000.0	my cat loves i	i started my c	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	0.14388489208	0.0	0.0	0.0	1.4388489208633095	
5.0	1315008000.0	great for fat c	we have three	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	0.7194244604	0.0	1.941747572f	0.0	2.6611720332471887	
2.0	1310515200.0	nearly killed t	too much of e	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	0.7194244604	0.7633587786	0.9708737864	0.0	2.453657025465376	
1.0	1220227200.0	changed form a	s with canide	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	[0, 0, 0, 0, 0, 0]	0.7633587786	1.941747572f	1.063829787f	0.0	3.7689361386755307	

Choix de l'algorithme de machine Learning

Un des premiers algorithmes que nous avons appliqués était celui du KNN. Pour l'appliquer, après avoir fait une prédiction des étiquettes selon de simples étiquettes « positif », « neutre » et « négatif », nous avons essayé d'appliquer l'algorithme de Machine Learning afin d'essayer de comprendre si on pouvait l'appliquer dans le cadre de notre analyse. Malheureusement, les résultats n'ont pas été concluants avec notamment un score d'accuracy maximal de 56 %. Nous n'avons ainsi pas retenu cet algorithme.

D'autre part, nous avons réalisé un modèle qui utilise du one hot encoding à partir d'une liste des 5 000 les plus présents dans tous les commentaires afin de prédire à l'aide d'un algorithme de KNN et de SVC une note sur 5.

Enfin, nous avons réalisé un modèle de machine Learning qui utilise une régression linéaire afin de déterminer le pourcentage de chaque sentiment présent dans les commentaires. Pour cela, nous avons fait du one hot encoding à partir de listes de mots que nous avons implémenté nous même pour chaque sentiment.

Entraînement du modèle

Concernant le modèle de la forêt aléatoire, afin d'évaluer notre modèle nous avons divisé notre base de données, en effet nous avons importé les colonnes Text et Sumaury que nous avons concaténer. À la suite de cela nous avons créé une fonction qui affecte un sentiment à chaque commentaire de la base de données, cette fonction a été réalisée grâce à la bibliothèque TextBlob de python. Pour finir nous avons entraîné un algorithme de forêt aléatoire, celui-ci est capable de lire et d'analyser les commentaires ainsi que les sentiments qui leur sont attribués. Grâce à cela l'algorithme est capable d'établir une prédiction.

Concernant les modèles de régression linéaire, nous avons aussi utilisé les colonnes Text et Sumaury que nous avons concaténer ainsi que les colonnes Target (liste en fonction du sentiment recherché) créées grâce au one hot encoding. Nous avons ensuite séparé ces données en données d'entraînement et de test (70% d'entraînement et 30% de test) afin d'entraîner notre modèle et de vérifier ses performances. Nous avons bien un modèle de régression linéaire par sentiment.

Résultats et Discussion

Développement

Afin de pouvoir présenter notre produit nous avons créé un site web permettant de retracer le parcours de notre projet ainsi que d'afficher les résultats de nos algorithmes. En voici un aperçu :



Notre site web est composé de cinq pages principales, chacune jouant un rôle spécifique pour présenter notre produit de manière complète et engageante.

Page d'accueil (Home) : La page d'accueil constitue l'interface principale de notre site. Elle offre une vue d'ensemble de notre produit, mettant en avant ses fonctionnalités principales et attirant l'attention des visiteurs dès leur arrivée. Nous utilisons un design attrayant et convivial pour encourager les utilisateurs à explorer davantage notre site.

Page À propos (About Us) : La page "À propos" est dédiée à retracer le parcours de nos membres. Nous partageons des informations sur notre équipe, notre expertise et notre expérience. Cette section met en valeur nos compétences, notre engagement et notre motivation à fournir un produit de haute qualité.

Page Contact : La page de contact est un moyen essentiel pour les visiteurs de communiquer avec notre équipe. Nous fournissons un formulaire de contact convivial où les utilisateurs peuvent nous envoyer des messages, poser des questions ou demander des informations supplémentaires.

Page Projets (Projects) : La page "Projets" est dédiée à la présentation de nos réalisations et de nos algorithmes. Nous mettons en valeur les projets sur lesquels nous avons travaillé, en fournissant des détails sur les problèmes que nous avons résolus, les méthodes que nous avons utilisées et les résultats obtenus.

Application

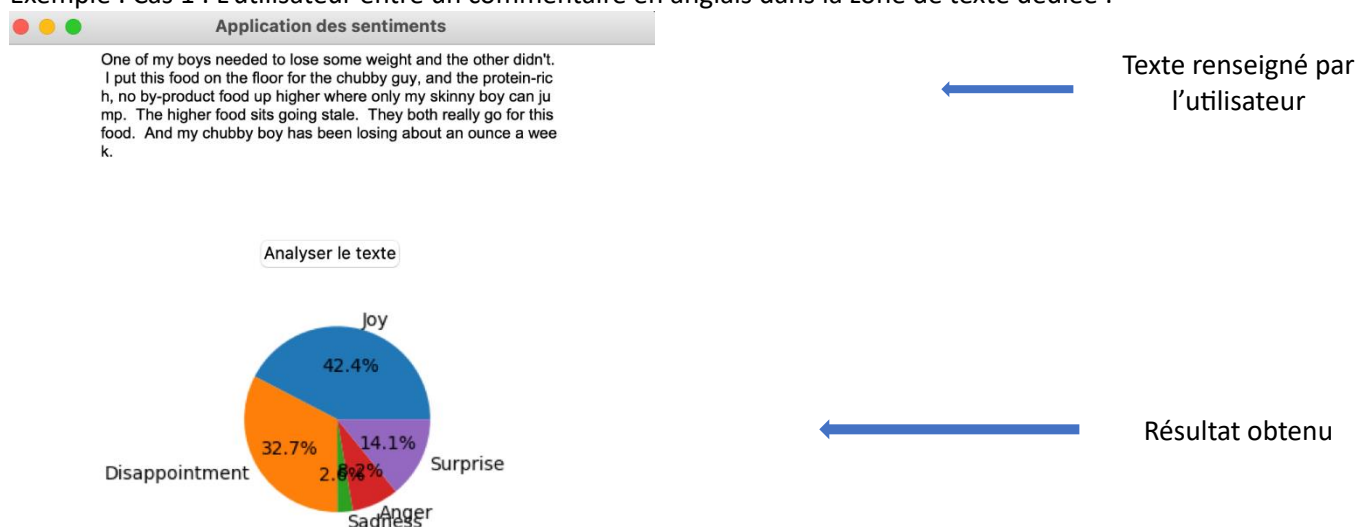
Notre application permet donc à l'utilisateur au choix :

- De rentrer une base de données ainsi que la colonne qu'il souhaite analyser
- De rentrer du texte (un commentaire)

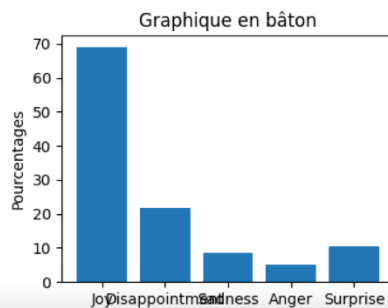
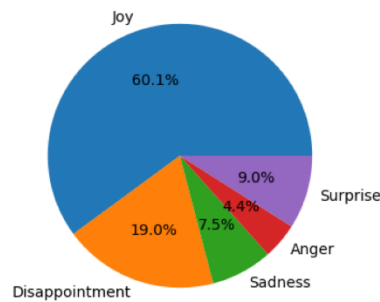
Une fois le texte où la base de données renseigné, notre programme traite les données afin de les rendre exploitable par nos modèles et réalise l'analyse des sentiments.

L'utilisateur obtient ainsi des graphiques qui rendent compte en pourcentage des sentiments exprimés par les commentaires (moyenne des pourcentages par sentiment si base de données renseignée). Si l'utilisateur renseigne un fichier à analyser, alors il est pertinent d'afficher les nuages de mots les plus fréquemment retrouvés dans les commentaires afin d'avoir une analyse plus approfondie.

Exemple : Cas 1 : L'utilisateur entre un commentaire en anglais dans la zone de texte dédiée :



Cas 2 : L'utilisateur renseigne le nom de son fichier (avec l'extension) ainsi que la colonne à analyser comme commentaire.

[illegible]

La base de données est aussi modifiée :

Summary	Text	prediction_joy	prediction_disappointment	prediction_sadness	prediction_anger	prediction_surprise
400: Good Quality Dog Food	I have bought several of the Vitality canned dog food 100.0	0.0	0.0	0.0	0.0	0.0
000: Not as Advertised	Product arrived labeled as Jumbo Salted Peanuts...tt 0.0	100.0	0.0	0.0	0.0	0.0
600: "Delight" says it all	This is a confection that has been around a few cent 59.08	45.54	3.69	11.38	19.69	19.69
200: Cough Medicine	If you are looking for the secret ingredient in Robitru: 59.08	45.54	3.69	11.38	19.69	19.69
600: Great Taffy	Great taffy at a great price. There was a wide assoi 100.0	0.0	0.0	0.0	0.0	0.0
200: Nice Taffy	I got a wild hair for taffy and ordered this five poun 100.0	0.0	0.0	0.0	0.0	0.0
400: Great! Just as good as the expensive brands!	This saltwater taffy had great flavors and was very s 100.0	0.0	0.0	0.0	0.0	0.0
200: Wonderful, tasty taffy	This taffy is so good. It is very soft and chewy. The 59.08	45.54	3.69	11.38	19.69	19.69
400: Yay Barley	Right now I'm mostly just sprouting this so my cats : 100.0	0.0	0.0	0.0	0.0	0.0
600: Healthy Dog Food	This is a very healthy dog food. Good for their digest 100.0	0.0	0.0	0.0	0.0	0.0
800: The Best Hot Sauce in the World	I don't know if it's the cactus or the tequila or ju 70.17	0.0	0.0	0.0	29.83	29.83
200: My cats LOVE this "diet" food better than their regular food	One of my boys needed to lose some weight and the 59.08	45.54	3.69	11.38	19.69	19.69

Exemple du nuage de point qui s'affiche :

Ce nuage représente les mots les plus utilisés en lien avec le sentiment de « joie ». Ainsi, l'utilisateur peut observer les éléments qu'il devrait conserver. Les mots non représentatifs tels que les pronoms et les adjectifs ont été supprimé.



Performance de l'algorithme

Pour évaluer des modèles, il y a plusieurs outils fournis par python pour comprendre les performances du modèle et son efficacité. Nous pouvons par exemple utiliser le coefficient R^2 . Ce coefficient permet de comprendre la précision du modèle. Il est compris entre 0 et 1.

D'autre part, la précision du modèle peut être établie par le calcul de l'accuracy.

L'évaluation de la forêt aléatoire est assez intuitive, en effet le terminal de l'IDE nous affiche les informations concernant la précision de notre algorithme :

```
In [2]: runcell('Forêts aléatoires', '/Users/anas/Desktop/MasterCamp/Projet_MasterCamp.py')
precision recall f1-score support
deçus      0.88      0.03      0.05      262
insatisfait 0.67      1.00      0.80     6073
neutre     0.87      0.01      0.03      929
non_disponible 0.43      0.04      0.07     1415
positif    0.85      1.00      0.92       33
accuracy   0.91      0.04      0.07     527
macro avg  0.67      0.35      0.32     9239
weighted avg 0.67      0.67      0.55     9239
```

Ici on observe que notre algorithme a une précision moyenne de 88%, ce qui est relativement fiable. Nous avons utilisé 20% des données pour les tests, et donc 80% des données sont utilisées pour l'entraînement du modèle.

Nous avons donc calculé les accuracy des modèles KNN et SVC qui déterminent le nombre d'étoile que possède un commentaire :

```
KNN accuracy: 0.5981383172132936
SVC accuracy: 0.656388626281725
```

Nous en avons déduit que le modèle SVC était plus adapté à notre projet car il a une meilleure précision. Cependant, nous avons finalement jugé pas très utile d'utiliser une prédiction d'une note de commentaire.

Concernant les 5 modèles de régression linéaires, nous avons calculé les R^2 afin de déterminer leur efficacité :

```
R2 Score modèle joie: 0.9997696718055805
R2 Score modèle déception: 0.996312019777414
R2 Score modèle tristesse: 0.9976701494750391
R2 Score modèle colère: 0.9754589499683444
R2 Score modèle surprise: 0.9971748895243675
```

Les modèles de régression linéaire ont très bon score R^2 et sont donc fiables et précis.

Comparaison avec d'autres approches existantes

Dans notre projet d'analyseur de sentiments, nous avons réalisé une comparaison approfondie avec d'autres approches existantes dans le domaine de l'analyse des sentiments. Cette comparaison nous permet de mettre en évidence les forces et les faiblesses de notre approche.

Les entreprises sur le marché se sont basées sur différentes méthodes pour obtenir un résultat similaire au nôtre, voici quelques-unes des méthodes déjà existantes sur le marché :

Certaines approches utilisent des règles prédéfinies pour classer les sentiments dans les textes. Cependant, ces approches sont souvent rigides et difficiles à généraliser, notre approche basée sur le Machine Learning, en revanche, permet au modèle de découvrir les motifs et les relations dans les données, offrant ainsi une meilleure capacité à traiter des textes complexes et à s'adapter à des contextes variés.

Les approches traditionnelles de l'analyse des sentiments reposent souvent sur des dictionnaires de mots pré-annotés avec des polarités positives, négatives ou neutres. Bien que ces approches soient simples à mettre en œuvre, elles sont limitées par leur dépendance aux dictionnaires, qui peuvent ne pas capturer les subtilités et les nuances du langage. Notre approche basée sur le Machine Learning associé aux dictionnaires offre une plus grande flexibilité et une meilleure capacité d'adaptation aux différentes nuances du langage.

Limitations et améliorations possibles

En ce qui concerne les limites de notre produit, nous avons identifié plusieurs points à prendre en considération, il est important de noter que notre produit n'est pas exempt de certaines contraintes qui pourraient avoir un impact sur son fonctionnement optimal.

Tout d'abord l'analyseur de sentiments peut être influencé par les variations linguistiques, tels que les expressions régionales, les argots, ou les tournures de phrases spécifiques à certaines communautés. Cela peut entraîner des erreurs d'interprétation des sentiments dans certains cas.

Les sentiments exprimés de manière sarcastique ou ironique peuvent être difficiles à détecter et à interpréter pour l'analyseur de sentiments. Ces formes d'expression peuvent conduire à des prédictions incorrectes ou à des classifications erronées.

En ce qui concerne les améliorations de notre produit, notre objectif est d'ajouter de nouvelles fonctionnalités pour permettre la compréhension de différentes langues. Actuellement, notre produit est uniquement disponible pour les entreprises anglophones, car ce marché présente un intérêt majeur. Toutefois, nous visons à élargir notre portée en développant des dictionnaires spécifiques à chaque pays, ce qui rendra notre produit fonctionnel dans le monde entier. Cette expansion linguistique nous permettra de répondre aux besoins des entreprises internationales et de saisir de nouvelles opportunités sur les marchés non anglophones.

Une autre amélioration potentielle consisterait à intégrer des informations contextuelles, telles que les informations sur l'auteur du texte, le contexte social ou les événements actuels, pour améliorer la précision de l'analyse des sentiments. Cette approche pourrait permettre de mieux comprendre les nuances et les variations des sentiments exprimés.

Conclusion

Contributions du projet

En ce qui concerne la contribution des membres de l'équipe au projet, nous avons divisé les tâches de manière à assurer une progression optimale. Rayen, Mordjane et Anas ont pris en charge le développement du back-end, où ils ont réalisé et entraîné les divers algorithmes de Machine Learning, tout en apportant leur contribution au rapport. De leur côté, Fodé et Gabriel ont pris en charge la partie

front-end, en réalisant le site web sur lequel notre projet sera présenté, tout en apportant également leur contribution au rapport. Felix à quant à lui été un intermédiaire entre le front-end et le back-end, il a travaillé sur les deux ainsi que sur le rapport.

Réflexion sur les résultats obtenus

Dans le cadre de notre projet, nous avons réalisé et de testé plusieurs algorithmes afin d'obtenir des résultats optimaux. Toutefois, nos premières tentatives n'ont pas produit des résultats suffisamment précis. Par exemple, nous avons obtenu une précision de seulement 56% avec l'algorithme KNN (K plus proches voisins) et de 88% avec l'algorithme de la forêt aléatoire.

Cependant, nous avons constaté que l'utilisation de la méthode du codage one-hot encoding s'est avérée plus efficace. En effet, cette approche a permis d'obtenir des résultats quasiment parfaits, avec une précision de 99%. Cette amélioration significative de la précision démontre l'importance de choisir la bonne méthode d'encodage des données pour obtenir des résultats de qualité dans notre projet d'analyse de sentiments.

Perspectives

Les perspectives d'amélioration de notre produit sont multiples, voici quelques-unes d'entre elles :

Dans notre projet actuel, nous avons principalement travaillé avec des textes en anglais. Cependant, une perspective intéressante serait d'élargir la portée de notre analyseur de sentiments pour prendre en charge des langues multiples. Cela nécessiterait de développer des dictionnaires spécifiques à chaque langue et de mettre en place des techniques de traduction automatique pour permettre une analyse de sentiments précise et adaptée à différents contextes linguistiques.

Une autre perspective intéressante serait d'intégrer notre algorithme d'analyse de sentiments dans des systèmes en temps réel. Cela pourrait être appliqué à des domaines tels que l'analyse des réseaux sociaux en temps réel, la surveillance de la réputation des marques ou la détection des tendances émergentes. Cette intégration permettrait d'obtenir des informations instantanées sur les sentiments des utilisateurs et d'agir en conséquence pour prendre des décisions stratégiques.

En conclusion, notre projet d'analyse de sentiments présente de nombreuses perspectives intéressantes, allant de l'amélioration de la précision à l'extension multilingue, en passant par l'intégration dans des systèmes en temps réel. Ces perspectives ouvrent des opportunités passionnantes pour continuer à développer et à améliorer notre algorithme, en le rendant plus puissant, polyvalent et utile pour une variété d'applications.