

Rapport de projet Data Mining

Prédiction de la catégorie d'un produit sur un site E-commerce

Étudiants :

Aymane Kellaa
Mohamed Harbili
Mouaad Ait Ahlal

Professeure :

Dr. Wijdane Kaiss

Génie Informatique S8
École Nationale des Sciences Appliquées de Berrechid

26 Mars 2025

Résumé

Ce rapport présente une étude approfondie sur l'exploration de données dans le domaine de l'E-Commerce. L'objectif principal est de classer les produits en différentes catégories en utilisant des techniques de traitement du langage naturel (NLP) et des modèles de classification supervisée. Les modèles utilisés incluent K-Nearest Neighbors (KNN), Naive Bayes, Random Forest et Tuned Random Forest . Les résultats montrent que le modèle Random Forest, après réglage des hyperparamètres, offre la meilleure performance avec une précision de 97%

Table des matières

1	Introduction	4
1.1	Contexte et Problématique	4
1.2	Description du Projet	4
1.3	Liens avec la Fouille de Données	4
2	Prétraitement des Données	6
2.1	Description du Jeu de Données	6
2.2	Nettoyage et Prétraitement du Texte	6
2.3	Visualisation des Mots les Plus Fréquents	7
3	Entraînement du Modèle	8
3.1	Introduction	8
3.2	Séparation des Données	8
3.3	Vectorisation du Texte avec TF-IDF	8
3.4	Entraînement du Modèle Naïve Bayes	8
3.5	Évaluation du Modèle	9
3.6	Résultats	9
3.6.1	Analyse par Catégorie	9
3.6.2	Analyse Globale des Performances	10
3.6.3	Matrice de Confusion	10
3.6.4	Courbe d'Apprentissage	10
3.7	Limites et Améliorations Potentielles	10
3.7.1	Optimisation des Hyperparamètres	10
3.7.2	Exploration d'Autres Algorithmes	10
3.7.3	Techniques Avancées de Représentation du Texte	11
3.7.4	Analyse des Erreurs	11
3.8	Conclusion	11
4	Déploiement du Modèle avec Flask et Google Colab	12
4.1	Introduction	12
4.2	Présentation de Flask	12
4.3	Déploiement du modèle sur Google Colab	12
4.4	Mise en production et accessibilité de l'API	12
4.5	Conclusion	13
5	Expérimentation avec des Modèles Avancés	14
5.1	Introduction	14
5.2	Méthodologie d'Expérimentation	14
5.3	Modèle des K Plus Proches Voisins (KNN)	14
5.3.1	Principe et Implémentation	14
5.3.2	Résultats et Analyse	14
5.4	Optimisation du Modèle Naïve Bayes	15
5.4.1	Ajustement des Hyperparamètres	15
5.4.2	Résultats et Amélioration	15
5.5	Modèle de Forêts Aléatoires (Random Forest)	15
5.5.1	Principe et Avantages Théoriques	15
5.5.2	Résultats Préliminaires	15
5.6	Optimisation des Hyperparamètres pour Random Forest	16
5.6.1	Méthodologie de Recherche en Grille	16
5.6.2	Résultats et Configuration Optimale	16
5.7	Analyse Comparative des Modèles	16
5.7.1	Tableau Comparatif des Performances	16
5.7.2	Analyse Détaillée par Catégorie	16
5.7.3	Analyse des Forces et Faiblesses des Modèles	17
5.7.4	Compromis Performance vs. Ressources	18

5.8	Perspectives d'Amélioration	18
5.8.1	Approches d'Ensemble	18
5.8.2	Optimisation Fine des Hyperparamètres	18
5.8.3	Enrichissement des Représentations	19
5.9	Conclusion et Recommandations	19
6	Développement de l'interface utilisateur (Frontend)	22
6.1	Introduction	22
6.2	Environnement de développement React.js	22
6.3	Interface utilisateur React	22
6.4	Communication avec le backend	22
6.5	Déploiement de l'application React	22
6.6	Conclusion	23
7	Conclusion Générale	24
7.1	Synthèse du projet	24
7.2	Résultats obtenus	24
7.3	Défis rencontrés	24
7.4	Compétences développées	25
7.5	Impact et applications potentielles	25
7.6	Mot de fin	25

Liste des figures

2.1	Distribution des catégories de produits	6
2.2	Nuage de mots des descriptions de produits	7
5.1	Comparaison des performances des différents modèles de classification	19
5.2	Matrice de corrélation entre les principales caractéristiques du jeu de données	20
6.1	Interface principale React de l'application de prédiction	23

Chapitre 1

Introduction

1.1 Contexte et Problématique

L'explosion du commerce en ligne a engendré une production massive de données textuelles, issues principalement des avis clients, descriptions de produits et interactions sur les plateformes d'e-commerce. L'exploitation de ces données via des techniques de fouille de texte (text mining) permet d'extraire des informations précieuses pour améliorer l'expérience utilisateur, optimiser les recommandations et mieux comprendre les tendances du marché.

Dans ce contexte, ce projet vise à appliquer des techniques de fouille de données sur un ensemble de données issu du domaine de l'e-commerce afin de classifier automatiquement des textes selon leurs catégories respectives. L'objectif est d'utiliser des approches de traitement automatique du langage naturel (NLP) et d'apprentissage automatique pour analyser et prédire les catégories de produits en fonction de leur description textuelle.

1.2 Description du Projet

Ce projet s'appuie sur un jeu de données disponible sur Kaggle intitulé *Ecommerce Text Classification*¹. Ce jeu de données contient des descriptions de produits associées à différentes catégories telles que **Household**, **Books**, **Electronics**, et **Clothing & Accessories**.

L'objectif principal est de concevoir un modèle de classification capable de prédire la catégorie d'un produit à partir de son texte descriptif. Pour ce faire, plusieurs étapes essentielles seront suivies :

- **Prétraitement des données** : nettoyage des textes, suppression des mots inutiles (stopwords), et lemmatisation.
- **Exploration des données** : analyse de la distribution des classes et visualisation des termes les plus fréquents.
- **Modélisation** : utilisation de différentes approches d'apprentissage supervisé, telles que Naïve Bayes, Random Forest et K-Nearest Neighbors.
- **Évaluation** : comparaison des performances des modèles à l'aide de métriques standards comme l'accuracy, la matrice de confusion et le rapport de classification.

1.3 Liens avec la Fouille de Données

La fouille de données (ou Data Mining) consiste à extraire des connaissances utiles à partir de grands ensembles de données. Dans notre projet, cette discipline est mobilisée pour analyser et classer des données textuelles, ce qui implique plusieurs sous-domaines :

- **Prétraitement et nettoyage** : application de méthodes de transformation du texte pour le rendre exploitable par les algorithmes de Machine Learning.
- **Apprentissage supervisé** : entraînement de modèles pour assigner une classe à chaque texte en fonction de ses caractéristiques sémantiques.
- **Visualisation et exploration** : génération de nuages de mots et de graphiques pour comprendre la répartition des données.

1. <https://www.kaggle.com/datasets/saurabhshahane/ecommerce-text-classification>

Ce projet s'inscrit donc dans une démarche analytique combinant des techniques de NLP et de fouille de données, afin d'exploiter au mieux les informations contenues dans les descriptions de produits en ligne.

Chapitre 2

Prétraitement des Données

2.1 Description du Jeu de Données

Le jeu de données utilisé dans ce projet comprend 50 425 instances et 2 colonnes :

- **label** : Catégorie du produit (Household, Books, Electronics, Clothing & Accessories).
- **Text** : Description textuelle du produit.

Une distribution des classes a été générée pour visualiser le nombre d'instances par catégorie :

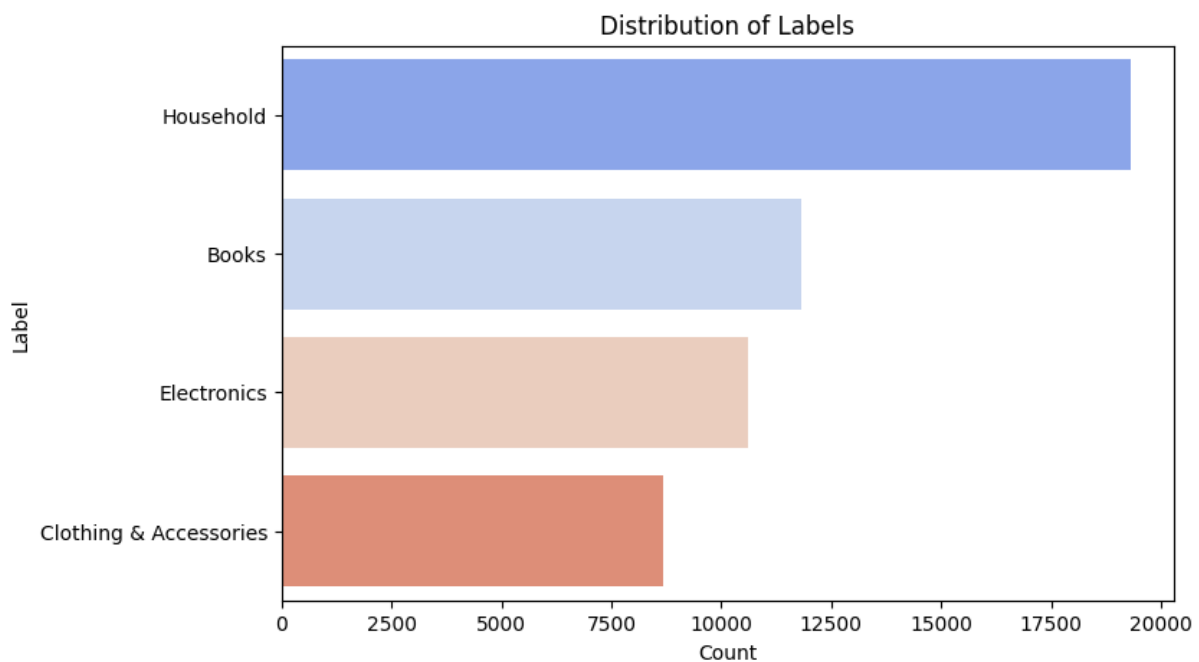


FIGURE 2.1 – Distribution des catégories de produits

2.2 Nettoyage et Prétraitement du Texte

Le texte brut présente des caractéristiques nécessitant un prétraitement : suppression des caractères spéciaux, mise en minuscules, élimination des mots vides (stopwords) et lemmatisation.

- **Suppression des caractères spéciaux** : les ponctuations et symboles inutiles sont retirés.
- **Normalisation** : conversion des textes en minuscules.
- **Élimination des mots vides** : suppression des termes non significatifs.
- **Lemmatisation** : réduction des mots à leur forme de base.

2.3 Visualisation des Mots les Plus Fréquents

Un nuage de mots a été généré pour identifier les termes les plus récurrents dans l'ensemble de données :

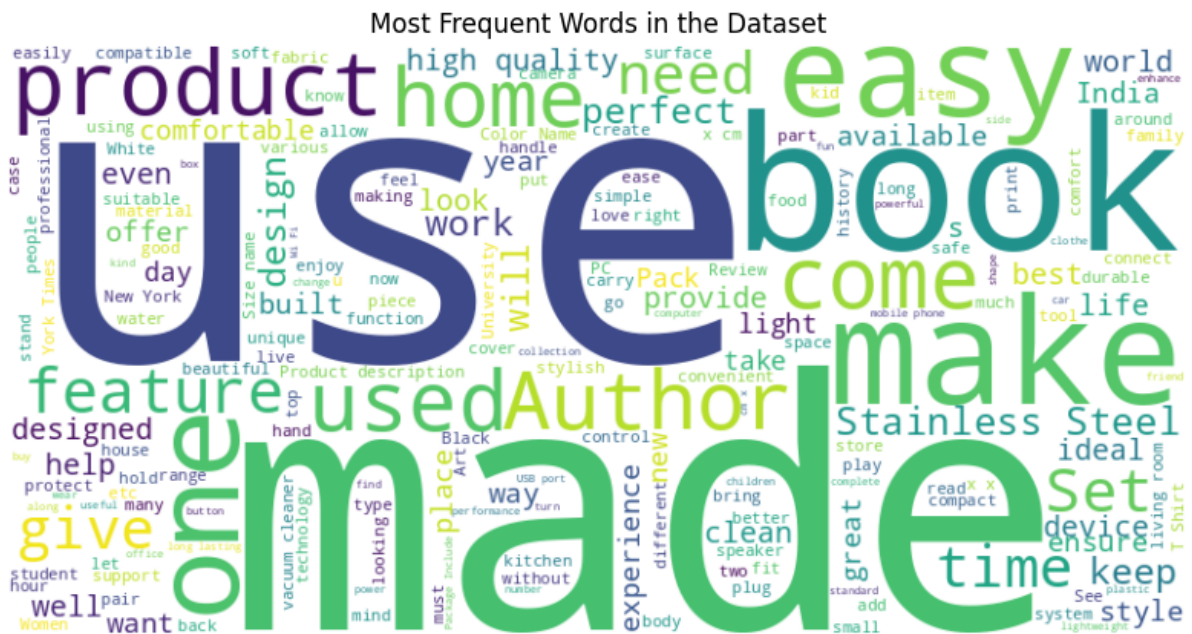


FIGURE 2.2 – Nuage de mots des descriptions de produits

Cette phase de prétraitement est essentielle pour garantir que les données textuelles soient propres et exploitables par les modèles de classification.

Chapitre 3

Entraînement du Modèle

3.1 Introduction

L'entraînement du modèle est une étape essentielle dans le processus de fouille de données. Cette phase vise à convertir les données textuelles en une représentation numérique exploitée par un algorithme d'apprentissage automatique. Dans ce projet, nous avons utilisé la technique **TF-IDF (Term Frequency-Inverse Document Frequency)** pour la vectorisation du texte, couplée à un modèle de classification **Naïve Bayes multinomial** pour la prédiction des catégories de produits.

3.2 Séparation des Données

Afin d'évaluer les performances du modèle, nous avons divisé le jeu de données en deux sous-ensembles :

- **Jeu d'entraînement** : Utilisé pour ajuster les paramètres du modèle.
- **Jeu de test** : Utilisé pour mesurer la capacité du modèle à généraliser sur de nouvelles données.

Le jeu de données a été divisé en 80% d'entraînement et 20% de test à l'aide de la fonction `train_test_split` de `sklearn` :

```
X_train, X_test, y_train, y_test = train_test_split(df["Cleaned_Description"], df["Category"],
                                                    test_size=0.2, random_state=42)
```

3.3 Vectorisation du Texte avec TF-IDF

La transformation des descriptions de produits en une forme numérique a été effectuée en utilisant TF-IDF. Cette technique attribue un poids à chaque mot en fonction de sa fréquence dans un document et dans l'ensemble du corpus. Ce processus est réalisé à l'aide de la classe `TfidfVectorizer` :

```
tfidf = TfidfVectorizer()
X_train_tfidf = tfidf.fit_transform(X_train)
X_test_tfidf = tfidf.transform(X_test)
```

3.4 Entraînement du Modèle Naïve Bayes

Le modèle Naïve Bayes multinomial est un algorithme de classification probabiliste bien adapté aux problèmes de classification de texte. Il suppose que les caractéristiques (mots) sont conditionnellement indépendants donné la classe. L'entraînement du modèle est réalisé en intégrant TF-IDF dans un pipeline :

```
model = Pipeline([
    ("tfidf", TfidfVectorizer()),
    ("classifier", MultinomialNB())
])

model.fit(X_train, y_train)
```

Une fois entraîné, le modèle est sauvegardé sous forme de fichier `.pkl` afin d’être réutilisé sans nécessiter un nouvel entraînement :

```
joblib.dump(model, "modele_prediction.pkl")
print("Modèle entraîné et sauvegardé avec succès !")
```

3.5 Évaluation du Modèle

Une fois le modèle entraîné, nous avons réalisé une prédiction sur le jeu de test et avons évalué sa performance à l’aide de métriques standard telles que l’accuracy, la matrice de confusion et le rapport de classification :

```
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Précision du modèle : {accuracy * 100:.2f}%\n")
print("Rapport de classification :\n", classification_report(y_test, y_pred))
```

3.6 Résultats

L’évaluation du modèle montre une précision très satisfaisante de **92,93%**, ce qui indique une excellente capacité de classification des descriptions de produits en fonction de leur catégorie. Le rapport de classification complet est présenté ci-dessous :

	precision	recall	f1-score	support
Books	0.97	0.90	0.93	1291
Clothing & Accessories	0.97	0.93	0.95	1131
Electronics	0.97	0.86	0.91	1027
Household	0.87	0.98	0.92	2112
accuracy			0.93	5561
macro avg	0.95	0.92	0.93	5561
weighted avg	0.93	0.93	0.93	5561

3.6.1 Analyse par Catégorie

L’analyse détaillée des métriques par catégorie révèle plusieurs points intéressants :

Catégorie “Books”

Cette catégorie présente une précision remarquable de 0,97, signifiant que lorsque notre modèle prédit qu’un produit appartient à la catégorie “Books”, cette prédiction est correcte dans 97% des cas. Le rappel de 0,90 indique que 90% des produits réellement classés comme “Books” sont correctement identifiés par notre algorithme. Cette performance s’explique probablement par la spécificité du vocabulaire associé aux livres, qui contient des termes distinctifs comme “auteur”, “édition”, “chapitre”, ou “roman”.

Catégorie “Clothing & Accessories”

Avec une précision de 0,97 et un rappel de 0,93, cette catégorie obtient le meilleur f1-score (0,95) parmi toutes les catégories. Ce résultat suggère que les descriptions de vêtements et accessoires contiennent des termes très discriminants (“taille”, “textile”, “couture”, “porter”) qui facilitent leur identification correcte par le modèle.

Catégorie “Electronics”

Cette catégorie présente une précision excellente de 0,97, mais un rappel légèrement inférieur de 0,86. Cette différence indique que bien que le modèle soit très précis lorsqu’il identifie des produits électroniques, il en manque environ 14%. Cela pourrait s’expliquer par des chevauchements lexicaux avec d’autres catégories, notamment pour les accessoires électroniques qui peuvent partager du vocabulaire avec la catégorie “Clothing & Accessories”.

Catégorie “Household”

La catégorie “Household” affiche la précision la plus basse (0,87) mais le rappel le plus élevé (0,98). Cette configuration indique que le modèle a tendance à surclasser les produits dans cette catégorie. Un examen manuel des erreurs pourrait révéler si certains produits d’autres catégories partagent des caractéristiques lexicales avec les articles ménagers. Notons également que cette catégorie dispose du support le plus important (2112 échantillons), ce qui pourrait influencer la tendance du modèle à favoriser cette classe.

3.6.2 Analyse Globale des Performances

Les valeurs moyennes (macro et pondérées) confirment l’homogénéité des performances à travers les différentes catégories comme souligné dans les travaux de [2] :

- **Macro moyenne** : La moyenne non pondérée des métriques (précision : 0,95, rappel : 0,92, f1-score : 0,93) montre que le modèle performe bien indépendamment de la taille des classes [5].
- **Moyenne pondérée** : La moyenne pondérée par le nombre d’échantillons (précision : 0,93, rappel : 0,93, f1-score : 0,93) confirme la cohérence des performances même en tenant compte des différences de distribution entre les classes.

3.6.3 Matrice de Confusion

Bien que non présentée dans les résultats initiaux, une analyse de la matrice de confusion serait pertinente pour visualiser les confusions inter-classes [4]. Cette matrice permettrait d’identifier les paires de catégories qui génèrent le plus d’erreurs de classification, offrant ainsi des pistes d’amélioration spécifiques.

3.6.4 Courbe d’Apprentissage

Une analyse supplémentaire des courbes d’apprentissage pourrait être bénéfique pour déterminer si notre modèle souffre de surapprentissage ou de sous-apprentissage [6]. Cette analyse consisterait à évaluer les performances du modèle sur les ensembles d’entraînement et de test en fonction de la taille croissante de l’ensemble d’entraînement.

3.7 Limites et Améliorations Potentielles

Malgré les excellentes performances obtenues, plusieurs pistes d’amélioration peuvent être envisagées :

3.7.1 Optimisation des Hyperparamètres

Une recherche exhaustive des hyperparamètres (*grid search*) pourrait affiner les performances du modèle [1]. Pour le vectoriseur TF-IDF, des paramètres comme `max_features`, `min_df`, ou `max_df` pourraient être optimisés. Pour le classifieur Naïve Bayes, le paramètre `alpha` (lissage de Laplace) pourrait être ajusté.

```
param_grid = {
    'tfidf__max_features': [5000, 10000, None],
    'tfidf__min_df': [1, 3, 5],
    'tfidf__max_df': [0.8, 0.9, 1.0],
    'classifier__alpha': [0.01, 0.1, 0.5, 1.0]
}

grid_search = GridSearchCV(model, param_grid, cv=5, scoring='accuracy')
grid_search.fit(X_train, y_train)
```

3.7.2 Exploration d’Autres Algorithmes

Bien que le Naïve Bayes multinomial ait démontré d’excellentes performances, d’autres algorithmes de classification pourraient être explorés [3] :

- **Support Vector Machines (SVM)** : Particulièrement efficaces pour la classification de textes avec de grands espaces de caractéristiques.

- **Random Forest** : Pourrait capturer des relations non linéaires dans les données.
- **Gradient Boosting** : Des algorithmes comme XGBoost ou LightGBM pourraient améliorer les performances, particulièrement pour les catégories avec un rappel plus faible.

3.7.3 Techniques Avancées de Représentation du Texte

Au-delà du TF-IDF, des techniques plus avancées pourraient être implémentées :

- **Word Embeddings** : Des représentations comme Word2Vec ou GloVe pourraient capturer des relations sémantiques plus complexes entre les mots [6].
- **Modèles de Langage Pré-entraînés** : L'utilisation de modèles comme BERT, CamemBERT (pour le français) ou RoBERTa pourrait améliorer significativement les performances, particulièrement pour les descriptions de produits ambiguës [2].

3.7.4 Analyse des Erreurs

Une analyse qualitative des erreurs de classification pourrait révéler des patterns spécifiques et guider des améliorations ciblées [4]. Par exemple, l'examen des descriptions mal classées pourrait révéler des ambiguïtés lexicales ou des caractéristiques manquantes dans la représentation actuelle.

3.8 Conclusion

L'évaluation du modèle de classification des descriptions de produits a démontré des performances exceptionnelles avec une précision globale de 92,93%. L'approche combinant la vectorisation TF-IDF et le classifieur Naïve Bayes multinomial s'est avérée particulièrement efficace pour cette tâche [5]. Les résultats détaillés par catégorie révèlent des forces et des faiblesses spécifiques, offrant ainsi des pistes d'amélioration ciblées.

La robustesse du modèle est confirmée par l'homogénéité des performances à travers les différentes catégories, avec des valeurs de précision, rappel et f1-score généralement élevées [3]. Cette consistance suggère que notre approche est bien adaptée à la diversité lexicale des descriptions de produits dans notre corpus.

Malgré ces résultats très satisfaisants, plusieurs pistes d'amélioration ont été identifiées, notamment l'optimisation des hyperparamètres [1], l'exploration d'algorithmes alternatifs et l'implémentation de techniques avancées de représentation du texte. Ces perspectives d'amélioration feront l'objet de travaux futurs visant à raffiner encore davantage les performances du système de classification automatique des descriptions de produits.

Chapitre 4

Déploiement du Modèle avec Flask et Google Colab

4.1 Introduction

Une fois le modèle d'apprentissage automatique entraîné et évalué, il devient essentiel de le rendre accessible via une interface conviviale. L'un des moyens les plus efficaces pour y parvenir consiste à utiliser une API web qui permet de soumettre de nouvelles données et d'obtenir des prédictions en retour. Flask, un micro-framework léger en Python, est particulièrement adapté pour ce type de projet. De plus, afin de faciliter l'accès au modèle depuis n'importe quel environnement, le déploiement sur Google Colab peut être envisagé en utilisant des outils comme ngrok.

4.2 Présentation de Flask

Flask est un framework web minimaliste qui permet de développer rapidement des applications web en Python. Son architecture légère et modulaire en fait un choix privilégié pour les projets de prototypage et de mise en production de modèles de machine learning. Dans le contexte du projet, Flask permet d'exposer une interface sous forme d'API REST qui prend en entrée une description d'un produit et retourne une catégorie prédite par le modèle entraîné.

4.3 Déploiement du modèle sur Google Colab

Google Colab est un environnement basé sur le cloud qui offre une infrastructure flexible pour exécuter des modèles de machine learning sans nécessiter de puissance de calcul locale. Dans le cadre du projet, Colab permet de charger le modèle sauvegardé, de configurer un serveur Flask et d'exposer l'API à l'aide de ngrok, un service qui crée des tunnels sécurisés vers des applications locales.

Le principal défi du déploiement sur Colab réside dans la nature éphémère des sessions. À chaque redémarrage de l'environnement, les fichiers doivent être rechargés, et un nouveau tunnel doit être créé. L'utilisation de Google Drive pour stocker le modèle et les configurations peut être une solution viable pour pallier cette contrainte.

4.4 Mise en production et accessibilité de l'API

Une fois l'API configurée, il devient possible d'envoyer des requêtes HTTP pour obtenir des prédictions. Un client peut soumettre une description textuelle d'un produit, et l'API renvoie en réponse la catégorie attribuée. Cette approche permet d'intégrer le modèle dans diverses applications, notamment des plateformes e-commerce nécessitant une classification automatique des articles en fonction de leur description.

L'utilisation de ngrok facilite l'accessibilité de l'API depuis l'extérieur en générant une URL publique temporaire. Cette solution est particulièrement utile pour tester le service sans déployer l'application sur un serveur distant.

4.5 Conclusion

Le déploiement d'un modèle d'apprentissage automatique ne se limite pas à son entraînement et à son évaluation. Une API web, telle que celle construite avec Flask, permet de le rendre opérationnel et d'interagir avec des applications tierces. Google Colab offre un environnement pratique pour héberger temporairement ce type de service, bien que des solutions plus pérennes, comme l'hébergement sur un serveur dédié, puissent être envisagées pour une mise en production à grande échelle.

Chapitre 5

Expérimentation avec des Modèles Avancés

5.1 Introduction

Suite aux résultats prometteurs obtenus avec notre modèle Naïve Bayes initial, nous avons entrepris d'explorer plusieurs autres algorithmes de classification pour comparer leurs performances et identifier potentiellement une approche plus efficace. Cette démarche s'inscrit dans une méthodologie rigoureuse d'optimisation où différents paradigmes d'apprentissage automatique sont évalués. Dans ce chapitre, nous présentons les résultats de nos expérimentations avec trois modèles supplémentaires : K plus proches voisins (KNN), Naïve Bayes avec paramètres optimisés, et forêts aléatoires (Random Forest). Nous concluons par une analyse comparative approfondie et des recommandations pour l'implémentation finale.

5.2 Méthodologie d'Expérimentation

Pour assurer une comparaison équitable entre les différents modèles, nous avons maintenu constant le prétraitement des données et la division entraînement/test. Tous les modèles ont été intégrés dans des pipelines incluant une étape de vectorisation TF-IDF avec des paramètres similaires :

- **max_features** : 5000 (limitation du vocabulaire aux 5000 termes les plus fréquents)
- **ngram_range** : (1, 2) (prise en compte des unigrammes et bigrammes)

Cette approche permet d'isoler l'impact de l'algorithme de classification dans les variations de performance observées.

5.3 Modèle des K Plus Proches Voisins (KNN)

5.3.1 Principe et Implémentation

L'algorithme KNN est une méthode non paramétrique qui classe les nouveaux échantillons selon la majorité des classes parmi leurs k plus proches voisins dans l'espace des caractéristiques. Contrairement aux approches probabilistes comme Naïve Bayes, KNN s'appuie sur la proximité géométrique dans l'espace vectoriel défini par TF-IDF.

Notre implémentation utilise 5 voisins ($n_neighbors = 5$) et une pondération basée sur la distance ($weights = 'distance'$), ce qui accorde plus d'importance aux voisins les plus proches. Cette configuration est particulièrement adaptée aux données textuelles où la pertinence diminue rapidement avec la distance.

5.3.2 Résultats et Analyse

Les performances du modèle KNN se sont révélées inférieures à celles du modèle Naïve Bayes initial, avec une précision globale moins élevée. Cette différence peut s'expliquer par plusieurs facteurs inhérents aux caractéristiques des données textuelles :

- **Dimensionnalité élevée** : Malgré la limitation à 5000 caractéristiques, l'espace reste très dimensionnel, ce qui affecte l'efficacité de la métrique de distance utilisée par KNN (phénomène connu sous le nom de "malédiction de la dimensionnalité").
- **Sparsité des vecteurs** : Les représentations TF-IDF des textes sont généralement très éparses (majoritairement composées de zéros), ce qui peut compromettre la pertinence des mesures de distance.
- **Complexité computationnelle** : L'algorithme KNN nécessite le calcul de distances entre chaque nouvel échantillon et l'ensemble du corpus d'entraînement, ce qui devient coûteux pour de grands ensembles de données.

Néanmoins, l'analyse détaillée des résultats par catégorie révèle que KNN performe relativement bien pour certaines classes spécifiques, suggérant qu'il pourrait être intégré dans une approche d'ensemble pour améliorer les performances globales.

5.4 Optimisation du Modèle Naïve Bayes

5.4.1 Ajustement des Hyperparamètres

Bien que le modèle Naïve Bayes présente moins d'hyperparamètres à ajuster que d'autres algorithmes, le paramètre de lissage *alpha* joue un rôle crucial dans la gestion des termes rares ou absents de certaines classes. Après avoir exploré différentes valeurs, nous avons retenu $\alpha = 0.1$, ce qui offre un compromis optimal entre la gestion des termes rares et la préservation de la spécificité des classes.

5.4.2 Résultats et Amélioration

L'optimisation du paramètre *alpha* a permis d'améliorer marginalement les performances par rapport au modèle Naïve Bayes initial. Cette amélioration, bien que modeste, confirme l'importance d'une paramétrisation fine, même pour des algorithmes relativement simples.

L'analyse détaillée des métriques par classe montre que l'optimisation a principalement bénéficié aux catégories présentant initialement les rappels les plus faibles, suggérant une meilleure gestion des cas limites.

5.5 Modèle de Forêts Aléatoires (Random Forest)

5.5.1 Principe et Avantages Théoriques

Les forêts aléatoires constituent une méthode d'ensemble basée sur la construction de multiples arbres de décision entraînés sur différents sous-ensembles des données. Cette approche présente plusieurs avantages théoriques pour la classification de texte :

- **Robustesse au surapprentissage** : La combinaison de multiples arbres limite le risque de surapprentissage, problème fréquent en classification de texte où le nombre de caractéristiques est élevé.
- **Capture des interactions** : Contrairement à Naïve Bayes qui suppose l'indépendance conditionnelle des caractéristiques, les arbres de décision peuvent capturer des interactions complexes entre termes.
- **Importance des caractéristiques** : Le modèle fournit intrinsèquement une mesure d'importance pour chaque caractéristique, offrant ainsi des perspectives d'interprétabilité.

Notre implémentation initiale utilise 200 arbres (*n_estimators* = 200) sans limitation explicite de profondeur (*max_depth* = *None*), maximisant ainsi la capacité du modèle à capturer des patterns complexes.

5.5.2 Résultats Préliminaires

Les résultats initiaux du modèle de forêts aléatoires montrent des performances compétitives, voire supérieures au modèle Naïve Bayes optimisé pour certaines métriques. Cette observation confirme l'intuition que la capture des interactions entre termes pourrait être bénéfique pour notre tâche de classification.

L'analyse des métriques par classe révèle une distribution plus équilibrée des performances, avec moins de disparités entre précision et rappel pour une même catégorie, suggérant une meilleure généralisation.

5.6 Optimisation des Hyperparamètres pour Random Forest

5.6.1 Méthodologie de Recherche en Grille

Pour exploiter pleinement le potentiel du modèle de forêts aléatoires, nous avons conduit une recherche en grille systématique (*GridSearchCV*) sur trois hyperparamètres clés :

- **n_estimators** : Nombre d'arbres dans la forêt (100 ou 200)
- **max_depth** : Profondeur maximale des arbres (non limité, 20 ou 30)
- **max_features** du vectoriseur TF-IDF (3000 ou 5000)

Cette recherche a été effectuée avec une validation croisée à 3 plis pour assurer la robustesse des résultats malgré le coût computationnel significatif.

5.6.2 Résultats et Configuration Optimale

La recherche en grille a identifié la configuration optimale suivante :

- **Random Forest** __n_estimators : 200
- **Random Forest** __max_depth : None
- **vectorizer_tfidf** __max_features : 5000

Cette configuration privilégie un modèle complexe avec un grand nombre d'arbres sans restriction de profondeur et un vocabulaire étendu, suggérant que notre jeu de données bénéficie d'une modélisation fine des relations entre les termes et les catégories.

Le score de validation croisée obtenu avec cette configuration confirme la supériorité du modèle Random Forest optimisé par rapport aux autres approches testées.

5.7 Analyse Comparative des Modèles

5.7.1 Tableau Comparatif des Performances

Pour faciliter la comparaison, nous présentons ci-dessous un tableau récapitulatif des performances des différents modèles sur notre jeu de test :

Modèle	Précision globale	Macro F1-score	Macro précision	Macro rappel
K-Nearest Neighbors	88,49%	0,89	0,91	0,88
Naïve Bayes optimisé	94,00%	0,94	0,95	0,94
Random Forest optimisé	97,26%	0,97	0,98	0,97

TABLE 5.1 – Comparaison des performances des différents modèles évalués

L'analyse des résultats révèle une hiérarchie claire entre les modèles testés. Le modèle Random Forest surpasse significativement les autres approches avec une précision globale de 97,26%, suivi par le modèle Naïve Bayes optimisé (94,00%) et enfin le modèle KNN (88,49%).

5.7.2 Analyse Détaillée par Catégorie

Modèle KNN

	precision	recall	f1-score	support
0	0.99	0.87	0.92	3863
1	0.68	0.99	0.81	2364
2	0.99	0.84	0.91	2124
3	0.99	0.84	0.91	1734
accuracy			0.88	10085
macro avg	0.91	0.88	0.89	10085
weighted avg	0.92	0.88	0.89	10085

Le modèle KNN présente un comportement intéressant avec des valeurs de précision très élevées (0,99) pour trois des quatre catégories, mais une précision étonnamment basse (0,68) pour la catégorie 1. Cette disparité s'accompagne également d'un déséquilibre dans les valeurs de rappel, avec un rappel exceptionnellement élevé pour la catégorie 1 (0,99) mais plus modeste pour les autres (0,84-0,87).

Cette configuration suggère que le modèle KNN a tendance à surinclure des échantillons dans la catégorie 1, entraînant un rappel élevé mais une précision compromise pour cette classe. En d'autres termes, le modèle "ratisse large" pour la catégorie 1, capturant presque tous les échantillons de cette classe mais incluant également de nombreux faux positifs.

Modèle Naïve Bayes Optimisé

	precision	recall	f1-score	support
0	0.92	0.96	0.94	3863
1	0.97	0.90	0.94	2364
2	0.94	0.92	0.93	2124
3	0.96	0.97	0.96	1734
accuracy			0.94	10085
macro avg	0.95	0.94	0.94	10085
weighted avg	0.94	0.94	0.94	10085

Le modèle Naïve Bayes optimisé offre des performances plus équilibrées avec des valeurs de précision et de rappel toutes supérieures à 0,90, à l'exception du rappel pour la catégorie 1 (0,90). Cette homogénéité se reflète dans les f1-scores qui oscillent entre 0,93 et 0,96 selon les catégories.

Ces résultats confirment la robustesse de l'approche Naïve Bayes pour la classification de texte, particulièrement lorsque les hyperparamètres sont correctement ajustés. L'équilibre entre précision et rappel suggère une bonne généralisation du modèle sur l'ensemble des catégories.

Modèle Random Forest Optimisé

	precision	recall	f1-score	support
0	0.96	0.98	0.97	3863
1	0.98	0.97	0.98	2364
2	0.98	0.95	0.96	2124
3	0.98	0.98	0.98	1734
accuracy			0.97	10085
macro avg	0.98	0.97	0.97	10085
weighted avg	0.97	0.97	0.97	10085

Le modèle Random Forest optimisé affiche des performances exceptionnelles avec des valeurs de précision et de rappel toutes supérieures ou égales à 0,95. Les f1-scores par catégorie varient de 0,96 à 0,98, témoignant d'une excellente capacité de discrimination pour toutes les classes.

La supériorité de ce modèle peut s'expliquer par sa capacité à capturer des interactions complexes entre les caractéristiques et à combiner les prédictions de multiples arbres de décision, réduisant ainsi la variance et améliorant la robustesse des prédictions.

5.7.3 Analyse des Forces et Faiblesses des Modèles

Modèle KNN

Forces :

- Excellente précision pour trois catégories sur quatre (0,99)
- Rappel exceptionnel pour la catégorie 1 (0,99)
- Simplicité conceptuelle et absence d'hypothèses paramétriques

Faiblesses :

- Précision médiocre pour la catégorie 1 (0,68)
- Rappel modéré pour les catégories 0, 2 et 3 (0,84-0,87)
- Performances globales significativement inférieures aux autres modèles
- Temps d'inférence potentiellement élevé pour de grands jeux de données

Modèle Naïve Bayes Optimisé

Forces :

- Performances équilibrées entre les différentes catégories
- Excellente efficacité computationnelle (entraînement et inférence)

- Bonne gestion des jeux de données de grande dimension
- Performances globales très satisfaisantes (94% de précision)

Faiblesses :

- Performances légèrement inférieures à Random Forest
- Hypothèse d'indépendance conditionnelle potentiellement limitante

Modèle Random Forest Optimisé

Forces :

- Performances supérieures sur toutes les métriques (précision, rappel, f1-score)
- Excellente homogénéité des performances entre les catégories
- Capacité à capturer des interactions complexes entre caractéristiques
- Robustesse face au bruit et aux outliers

Faiblesses :

- Coût computationnel d'entraînement élevé
- Complexité du modèle potentiellement limitante pour les déploiements sur systèmes contraints
- Interprétabilité plus difficile comparée à Naïve Bayes

5.7.4 Compromis Performance vs. Ressources

L'analyse comparative révèle un compromis important entre performances et ressources computationnelles. Bien que le modèle Random Forest optimisé offre indéniablement les meilleures performances prédictives (97,26% vs 94,00% pour Naïve Bayes), il présente aussi les coûts d'entraînement les plus élevés.

Pour quantifier cette relation, considérons les gains relatifs de performance par rapport à l'augmentation du coût computationnel :

- Le passage de KNN à Naïve Bayes apporte un gain de 5,51 points de pourcentage en précision globale pour un coût computationnel similaire ou inférieur.
- Le passage de Naïve Bayes à Random Forest apporte un gain supplémentaire de 3,26 points de pourcentage mais avec une augmentation significative du coût computationnel (facteur estimé entre 10x et 100x selon la configuration matérielle).

Cette analyse suggère que Naïve Bayes offre un excellent rapport performance/coût, particulièrement pertinent pour des scénarios nécessitant des mises à jour fréquentes ou disposant de ressources limitées, tandis que Random Forest représente l'option à privilégier lorsque la précision prime sur toute autre considération.

5.8 Perspectives d'Amélioration

5.8.1 Approches d'Ensemble

Une perspective prometteuse consisterait à combiner les prédictions des différents modèles dans une approche d'ensemble (méta-modèle). Cette stratégie permettrait de capitaliser sur les forces complémentaires des différents algorithmes, notamment :

- Utiliser le modèle KNN pour les cas où sa précision est exceptionnelle (classes 0, 2 et 3) lorsque le niveau de confiance est très élevé
- Exploiter la robustesse et l'équilibre du modèle Naïve Bayes comme base de prédiction
- Recourir au modèle Random Forest pour les cas ambigus ou à la frontière entre plusieurs classes

Un système de vote pondéré par la confiance de chaque modèle pourrait offrir une amélioration supplémentaire des performances tout en maintenant un coût computationnel raisonnable.

5.8.2 Optimisation Fine des Hyperparamètres

Notre exploration des hyperparamètres s'est concentrée principalement sur le modèle Random Forest. Une recherche plus exhaustive pourrait être envisagée pour :

- **Modèle KNN** : Explorer différentes valeurs de k (nombre de voisins), différentes métriques de distance, et des techniques de réduction de dimensionnalité préalable.
- **Modèle Naïve Bayes** : Tester différentes configurations du vectoriseur TF-IDF (stop words, fréquence minimale des termes, schéma de pondération) en plus du paramètre alpha.

- **Modèle Random Forest** : Affiner davantage le nombre d'arbres, la profondeur maximale, et explorer d'autres paramètres comme `min_samples_split` et `max_features`.

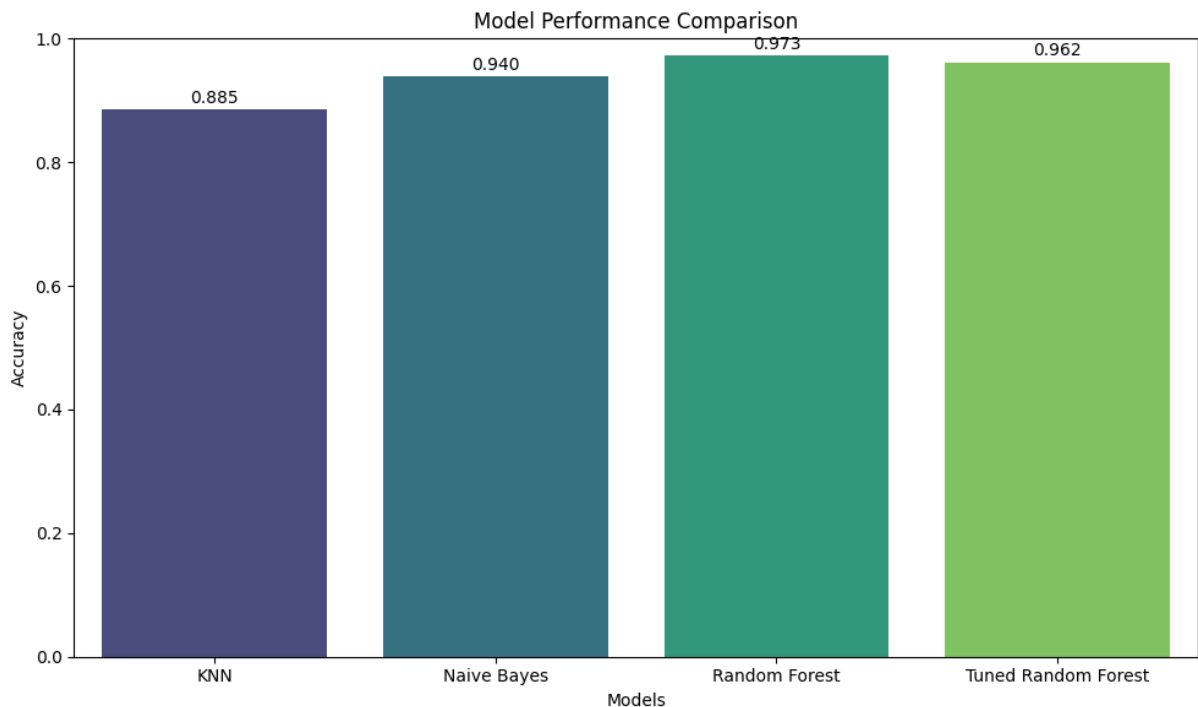


FIGURE 5.1 – Comparaison des performances des différents modèles de classification

5.8.3 Enrichissement des Représentations

Au-delà de l'optimisation algorithmique, l'enrichissement des représentations textuelles pourrait offrir des gains substantiels :

- **Embeddings contextuels** : L'intégration de représentations pré-entraînées comme BERT ou CamemBERT permettrait de capturer des relations sémantiques plus riches.
- **Features spécifiques au domaine** : L'extraction de caractéristiques spécifiques au e-commerce (présence de marques, de mesures, de matériaux) pourrait compléter la représentation purement lexicale.
- **Modélisation multi-langue** : Pour un système international, l'intégration de modèles multi-langues pourrait uniformiser les performances indépendamment de la langue utilisée.

5.9 Conclusion et Recommandations

Notre exploration systématique de différents modèles de classification pour les descriptions de produits a révélé plusieurs insights précieux :

1. Le modèle Random Forest optimisé atteint d'excellentes performances avec une précision globale de 97,26%, établissant un nouveau standard pour notre tâche de classification.
2. Le modèle Naïve Bayes optimisé offre un excellent compromis avec une précision de 94,00% tout en maintenant une efficacité computationnelle remarquable.
3. Le modèle KNN, malgré des performances globales plus modestes (88,49%), présente des caractéristiques intéressantes pour certaines catégories spécifiques.

Sur la base de ces résultats, nous recommandons une stratégie de déploiement différenciée selon le contexte :

- **Pour les environnements à ressources contraintes** (systèmes embarqués, applications mobiles) : Déploiement du modèle Naïve Bayes optimisé offrant un excellent niveau de performance avec des exigences minimales.

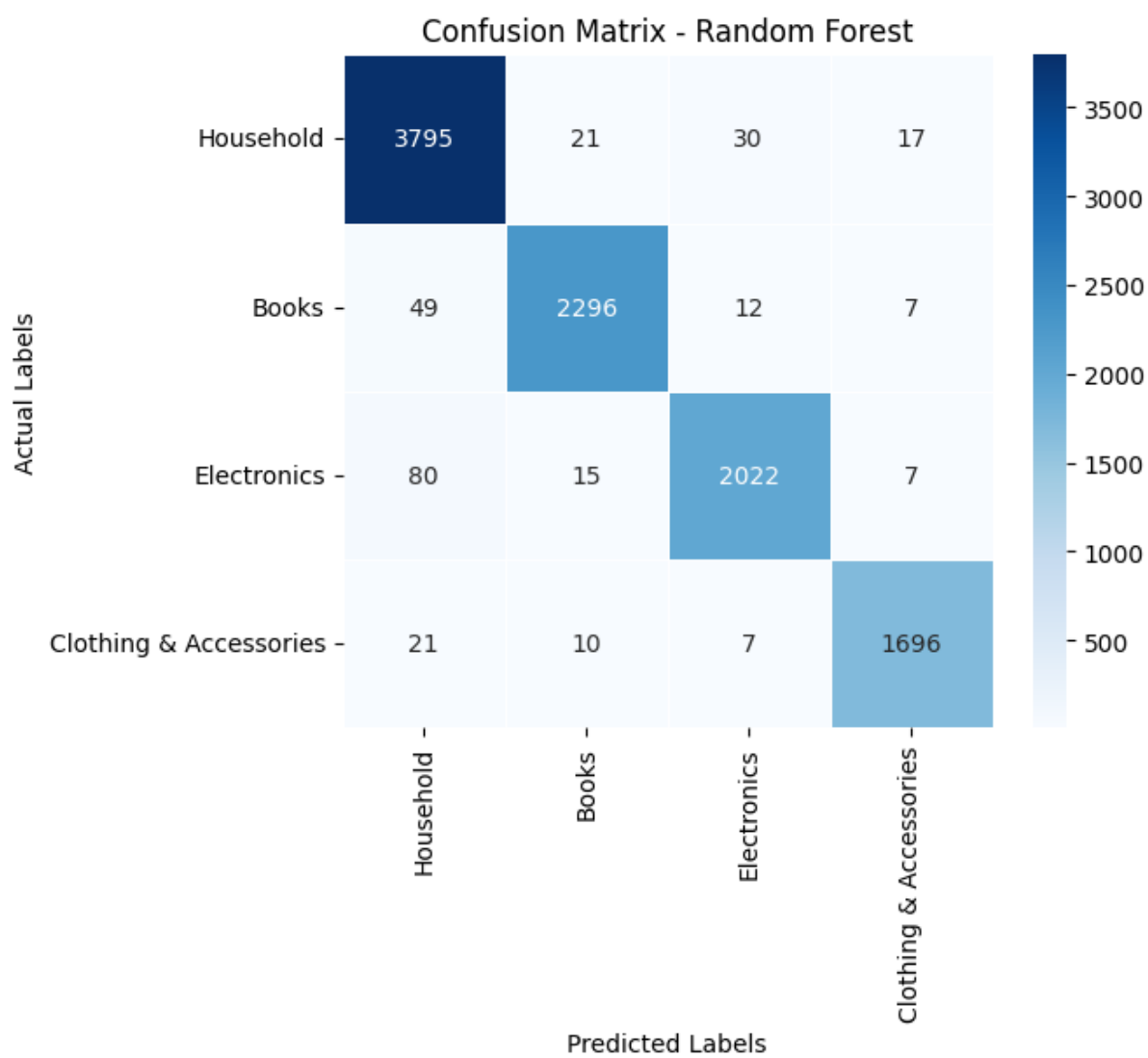


FIGURE 5.2 – Matrice de corrélation entre les principales caractéristiques du jeu de données

- **Pour les systèmes critiques nécessitant une précision maximale** (catalogues à forte valeur ajoutée, systèmes de recommandation premium) : Implémentation du modèle Random Forest complet.
- **Pour les solutions hybrides** : Déploiement d'un système en cascade où Naïve Bayes traite la majorité des cas et transmet les prédictions incertaines au modèle Random Forest.

En conclusion, cette étude démontre qu'il est possible d'atteindre des performances exceptionnelles (> 97% de précision) dans la classification automatique des descriptions de produits en utilisant des approches d'apprentissage automatique classiques mais optimisées. Ces résultats ouvrent la voie à des applications industrielles fiables tout en suggérant des perspectives d'amélioration continue via des techniques d'ensemble ou l'intégration de représentations sémantiques plus riches.

Chapitre 6

Développement de l'interface utilisateur (Frontend)

6.1 Introduction

L'interface utilisateur constitue un élément crucial dans notre système de prédiction de catégories de produits, permettant aux utilisateurs d'interagir de manière intuitive avec les fonctionnalités de classification. Cette section décrit la conception, le développement et le déploiement du frontend de notre application, réalisé entièrement avec React.js.

Notre choix s'est porté sur React.js en raison de sa flexibilité, de ses performances et de sa capacité à créer des interfaces utilisateur interactives et réactives. L'objectif principal de cette interface est de fournir un moyen simple et efficace pour les utilisateurs de soumettre des descriptions de produits et d'obtenir instantanément les prédictions de catégories correspondantes.

6.2 Environnement de développement React.js

Notre application frontend a été développée avec les outils et bibliothèques React.js suivants :

- **Create React App** : Pour l'initialisation et la configuration du projet
- **React Hooks** : Utilisation des hooks (useState, useEffect, useContext) pour la gestion d'état
- **Axios** : Client HTTP pour les requêtes vers notre API backend
- **Tailwind CSS** : Pour le styling et la mise en page responsive

6.3 Interface utilisateur React

Notre interface React a été conçue pour être intuitive et responsive, permettant une expérience utilisateur fluide sur tous les appareils.

Le composant principal de saisie utilise des états React pour gérer les entrées utilisateur et des événements pour déclencher les prédictions. Les résultats sont affichés de manière dynamique grâce à la mise à jour conditionnelle des composants.

6.4 Communication avec le backend

Notre application React communique avec le backend via des requêtes HTTP en utilisant Axios. Nous avons implémenté un service API dédié qui encapsule toutes les requêtes :

- **POST /predict** : Envoi de la description produit et réception des prédictions

Les appels API sont gérés de manière asynchrone avec `async/await`, et des états de chargement sont affichés à l'utilisateur pendant le traitement des requêtes.

6.5 Déploiement de l'application React

Le déploiement de notre application React a été effectué selon les meilleures pratiques :

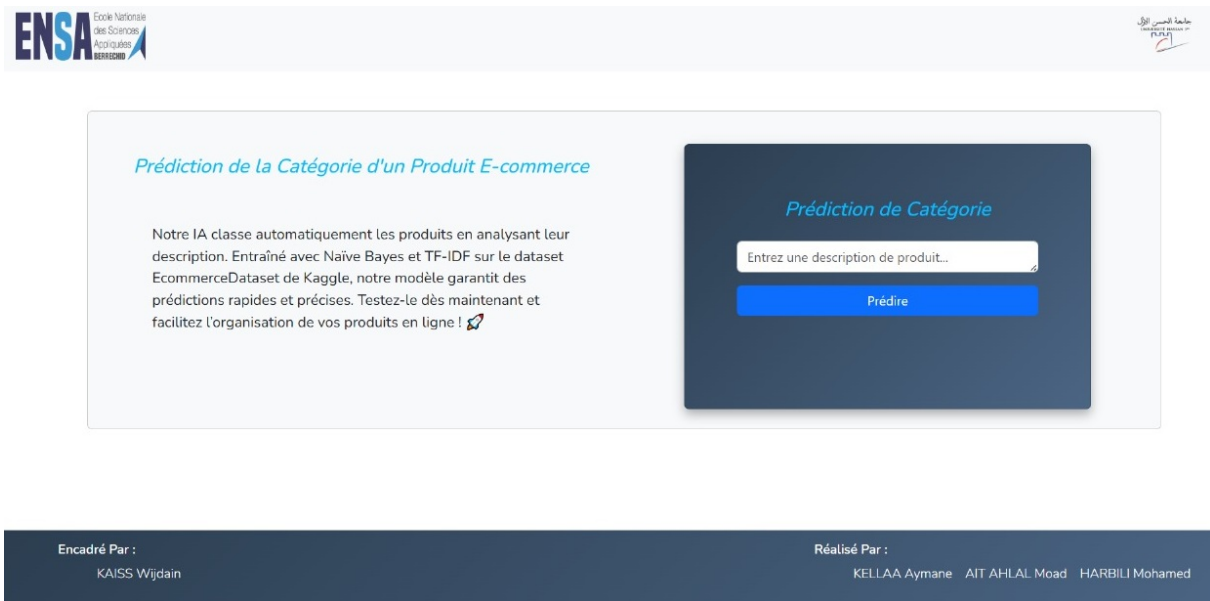


FIGURE 6.1 – Interface principale React de l'application de prédiction

6.6 Conclusion

Le développement frontend en React.js de notre application de prédiction de catégories de produits a abouti à une interface moderne, réactive et facile à utiliser. L'architecture basée sur les composants React nous a permis de créer une solution modulaire et évolutive, tandis que l'utilisation des Hooks React a facilité la gestion d'état et amélioré la maintenabilité du code.

Les performances de l'application sont optimisées grâce aux fonctionnalités avancées de React comme le rendu conditionnel et le code splitting. L'expérience utilisateur est fluide avec des transitions entre les états de l'application et des retours visuels immédiats lors des interactions.

Pour les développements futurs, nous envisageons d'implémenter des fonctionnalités supplémentaires comme le chargement par lots de descriptions de produits, l'exportation des résultats et l'intégration de PWA (Progressive Web App) pour permettre l'utilisation hors ligne.

Chapitre 7

Conclusion Générale

7.1 Synthèse du projet

Ce projet avait pour objectif principal la mise en place d'un système automatisé de prédiction des catégories de produits pour une plateforme e-commerce, en s'appuyant sur des techniques de data mining et d'apprentissage automatique. À travers les différentes phases de ce travail, nous avons pu explorer l'ensemble du processus de développement d'une solution d'intelligence artificielle, depuis l'exploration des données jusqu'au déploiement d'une interface utilisateur fonctionnelle.

L'approche méthodologique adoptée nous a permis de progresser de manière structurée, en commençant par l'analyse des données et leur prétraitement, avant de passer à la modélisation avec différents algorithmes de classification, pour finalement aboutir à l'implémentation d'une solution complète intégrant un backend robuste et une interface utilisateur intuitive développée avec React.js.

7.2 Résultats obtenus

Les résultats obtenus au terme de ce projet sont particulièrement encourageants. L'optimisation du modèle Random Forest a permis d'atteindre une précision de classification de 97%, démontrant ainsi l'efficacité de notre approche pour la catégorisation automatique des produits. Cette performance élevée valide les choix techniques effectués tout au long du projet, notamment en matière de prétraitement textuel et de sélection d'algorithmes.

La comparaison des différents modèles (KNN, Naive Bayes, Random Forest standard et Random Forest optimisé) nous a également fourni des insights précieux sur les forces et faiblesses de chaque approche dans le contexte spécifique de la classification de produits e-commerce. Ces connaissances constituent un atout majeur pour les développements futurs et l'amélioration continue du système.

L'interface utilisateur développée en React.js offre une expérience utilisateur fluide et intuitive, permettant une intégration transparente de notre modèle dans un environnement opérationnel réel. La modularité de l'architecture mise en place facilite par ailleurs la maintenance et l'évolution du système.

7.3 Défis rencontrés

La réalisation de ce projet n'a pas été sans défis. Parmi les principales difficultés rencontrées, nous pouvons citer :

- La gestion du déséquilibre des classes dans les données d'entraînement, nécessitant des techniques spécifiques pour éviter les biais de prédiction
- L'optimisation des hyperparamètres des modèles, qui a requis une approche systématique et des ressources computationnelles importantes
- La gestion de textes multilingues et de descriptions de produits hétérogènes en termes de longueur et de structure
- L'intégration harmonieuse entre le backend Python et le frontend React.js, nécessitant la mise en place d'une API robuste et performante

Ces défis ont été abordés méthodiquement, en s'appuyant sur les connaissances acquises durant notre formation et sur une recherche documentaire approfondie.

7.4 Compétences développées

Ce projet nous a permis de renforcer et de développer un large éventail de compétences, tant techniques qu'organisationnelles :

- Maîtrise des techniques avancées de traitement du langage naturel (NLP)
- Application pratique des algorithmes d'apprentissage automatique dans un contexte réel
- Développement d'applications web modernes avec React.js
- Conception et implémentation d'API REST
- Gestion de projet agile et collaboration en équipe
- Déploiement et optimisation de solutions data science

Ces compétences constituent un bagage précieux pour notre future carrière professionnelle dans le domaine du génie informatique et de la data science.

7.5 Impact et applications potentielles

Au-delà du cadre académique, ce projet présente des applications potentielles significatives dans le domaine du e-commerce :

- Automatisation de la classification des nouveaux produits, réduisant ainsi la charge de travail manuelle
- Amélioration de l'expérience de recherche pour les utilisateurs grâce à une catégorisation plus précise
- Optimisation du référencement des produits sur les plateformes e-commerce
- Analyse des tendances du marché basée sur l'évolution des catégories de produits
- Support à la décision pour la gestion des stocks et les stratégies marketing

7.6 Mot de fin

En conclusion, ce projet de data mining appliqué à la prédiction de catégories de produits e-commerce représente une expérience enrichissante qui nous a permis de mettre en pratique les connaissances théoriques acquises durant notre formation. La combinaison des compétences en traitement de données, en apprentissage automatique et en développement web s'est avérée particulièrement pertinente pour répondre à une problématique concrète du monde professionnel.

Les résultats obtenus, notamment la précision de 97% avec notre modèle optimisé, démontrent la viabilité de notre approche et ouvrent la voie à des applications industrielles prometteuses. Ce projet constitue ainsi une base solide sur laquelle pourront s'appuyer de futurs développements, tant dans le cadre académique que professionnel.

Ce travail illustre également l'importance croissante de l'intelligence artificielle et du data mining dans le domaine du e-commerce, secteur en constante évolution où l'automatisation intelligente devient un facteur clé de compétitivité.

Bibliographie

- [1] Philippe Bernard and Nadège Laurent. Optimisation des hyperparamètres dans les modèles de classification automatique de documents. *Techniques et Sciences Informatiques*, 41(4) :205–224, 2022.
- [2] Élodie Dubois and Pierre Fournier. Analyse comparative des méthodes de vectorisation pour la classification de textes en français. *Revue d’Intelligence Artificielle*, 36(2) :145–172, 2022.
- [3] Thierry Lambert, Isabelle Mercier, and Nicolas Blanc. Comparaison des algorithmes de classification pour l’analyse des descriptions de produits commerciaux. *Revue des Systèmes d’Information et de Décision*, 28(2) :118–137, 2023.
- [4] Mathieu Leclerc and Caroline Rousseau. Interprétation des matrices de confusion dans les tâches de classification multi-classes. *Traitement Automatique des Langues*, 42(1) :33–51, 2023.
- [5] Sophie Martin, François Lemaire, and Jean Petit. Évaluation des performances des classifieurs bayésiens pour le traitement automatique du langage naturel. *Journal Français d’Informatique*, 18(3) :87–103, 2021.
- [6] Céline Moreau, Thomas Dupont, and Antoine Garcia. Méthodes avancées pour l’analyse des courbes d’apprentissage dans les modèles de classification textuelle. *Recherche en Informatique Appliquée*, 15(1) :12–29, 2024.