

Classification supervisée et traitement du langage naturel

Durand Enzo[†], Serraf Dan[‡]

Abstract. Nous avons mené une campagne d'expérience afin de comprendre 2 jeux de données, puis de prédire l'appartenance à une classe de nouvelles données. La première étape est de comprendre les données textuelles grâce à plusieurs techniques de visualisation de données. Il faut ensuite faire beaucoup de différents tests, car les intuitions ne suffisent pas toujours. Pour cela, il faut réfléchir à l'organisation d'un pipeline car on ne peut pas tester toutes les combinaisons possibles de paramètres. Il faut veiller à utiliser les métriques les plus intéressantes, mais aussi faire attention au sur-apprentissage.

1. Jeu de données des présidents

Dans cette partie, on s'intéressera au jeu de données sur les présidents Jacques Chirac et François Mitterrand. Nous avons un fichier de données contenant des blocs de discours entre ces deux hommes, avec les étiquettes associées. Un deuxième fichier contient des blocs de discours mais cette fois-ci non étiquetés. Le but de ce projet est d'arriver à prédire qui est le locuteur de chacune des phrases contenu dans le fichier de test.

1.1. Analyse des données

Lorsque l'on travaille sur un problème de classification, une des premières étapes est de visualiser la répartition des classes dans nos données. En effet, cela est très important car un modèle de machine learning fera toujours au plus simple afin de maximiser un score de précision. Par exemple dans notre cas, nous avons 87% d'exemples pour Jacques Chirac et 13% d'exemples pour François Mitterrand. Cela veut dire que notre classifieur peut avoir 87% de précision en prédisant que toutes les phrases proviennent de la classe majoritaire.

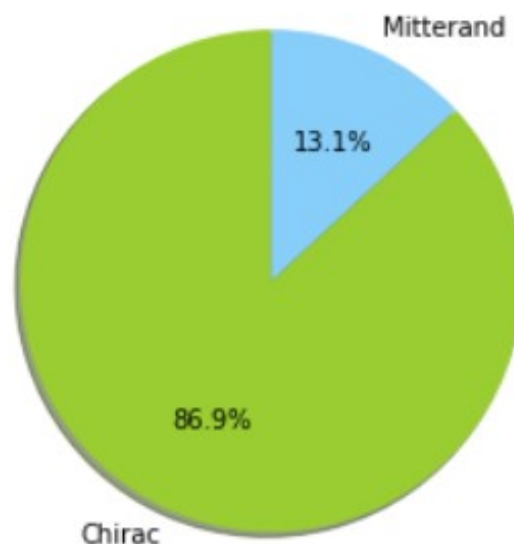
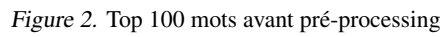


Figure 1. Répartition classes Présidents.



The chart displays the distribution of word frequencies. The x-axis represents the word index from 0 to 60,000, and the y-axis represents the logarithm of the number of occurrences per word, ranging from 0 to 10. The data shows a long tail of low-frequency words (log value around 0.5) for the first 25,000 words, followed by a gradual increase in frequency for the next 15,000 words (log values between 1 and 2), and finally a sharp spike to a log value of approximately 8.5 for the last word in the list.

Ces mots ne nous seront probablement pas utiles afin de différencier qui est le locuteur dans une phrase. Deux choix s'ouvrent donc à nous, utiliser des listes de **stopwords** prédéfinies, ou utiliser le paramètre **max_df** de notre `vectorizer`. La différence étant que le paramètre **max_df** permet de supprimer les stopwords propre au dataset.

Nous avons affiché plusieurs wordcloud mais même en supprimant des **stopwords**, il est compliqué de voir les différences dans les données. C'est pourquoi nous avons décidé d'utiliser le odds ratio afin d'afficher des wordclouds assez différents suivant les deux classes.

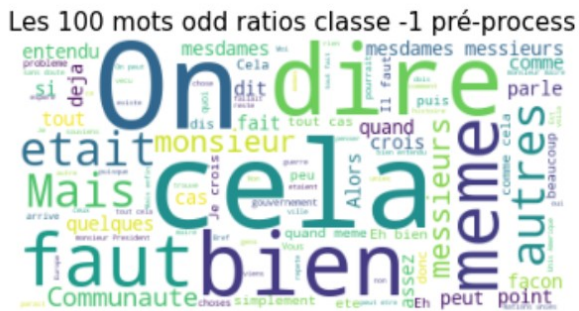


Figure 4. Top 100 mots après pré-processing Mitterrand.



Figure 5. Top 100 mots après pré-processing Chirac.

1.2. Recherche des paramètres optimaux pour la représentation des données

1.2.1. REPRÉSENTATION DES DONNÉES

Il y a plusieurs moyens de représenter les données textuelles afin d'utiliser des modèles de machine learning. Nous avons opté pour le `TfidfVectorizer` de `sklearn` après avoir testé plusieurs possibilités. Le `TfidfVectorizer` est en fait la combinaison d'un `CountVectorizer` et d'un `TfidfTransformer`. Il permet de représenter les mots sous un format statistique qui est intéressant pour la suite.

1.2.2. ORGANISATION

Nous avons mene une campagne d'expérience afin de trouver les paramètres optimaux pour notre problème. Premièrement, nous avons décidé d'utiliser les fonctions de gridsearch et pipeline de la librairie sklearn. Ces fonctions permettent d'effectuer une recherche de parametres optimaux suivant une métrique et une liste de parametres. Il est possible d'utiliser gridsearch sur des vectorizer (countvectorizer et tfidfvectorizer) ainsi que sur des modeles (LR, MultinomialNB et SVC). Il est possible de paralléliser les opérations grâce au paramètre **njobs** de gridsearch, ce qui est non negligeable car une des ressources limitantes est le temps.

Un des principaux problèmes est que le nombre de combinaisons de paramètres à tester augmente très rapidement. On ne peut donc pas faire un gridsearch sur l'ensemble des paramètres, il faut faire plusieurs choix.

Nous avons décidé de séparer notre recherche en 3 parties :

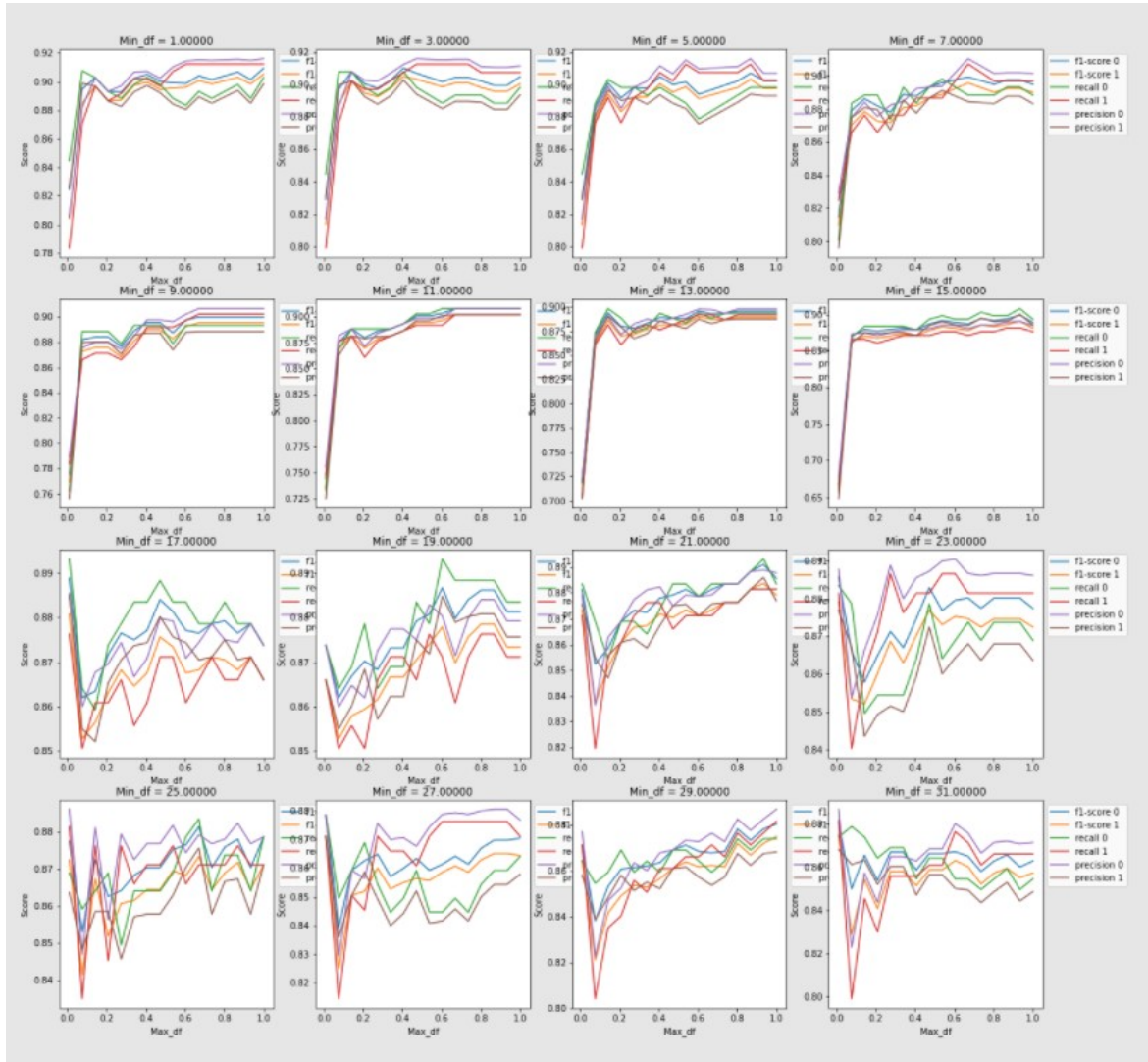
- La recherche des paramètres optimaux pour les données continues (**min_df**, **max_df**, **max_features**).
- La recherche des paramètres booléens ou discret (**use_idf**, **smooth_idf**, **sublinear_idf**, **lowercase**, **strip_accents**, **stopwords**, **ngram_range**).
- La recherche des paramètres des modèles de machine learning (**C**, **alpha**, **class_weights**, **max_iter**).

Afin d'éviter de faire un gridsearch comparant des modèles en sur-apprentissage, il faut prendre quelques

précautions. Lors du gridsearch sur les vectorizers nous utiliserons un modèle prédéfinis dans le pipeline. Ce modèle sera un SVM linéaire, ayant comme paramètres **class_weights="balanced"** et **C=100**. On utilisera aussi un **max_iter=1000** afin de gagner du temps sur la recherche de paramètre. La métrique la plus intéressante est le **f1** score lorsque nous avons des classes non-équilibrées. Nous avons donc utilisé cette métrique dans notre gridsearch. Les métriques **f1_micro** et **rog_auc** ont aussi été essayés mais elles donnent des résultats moins concluants.

1.2.3. PARAMÈTRES CONTINUS

Après avoir effectué quelques expériences, il semble que les paramètres continus (**min_df**, **max_df**, **max_features**) soient ceux qui influent le plus sur nos résultats. C'est pourquoi nous décidons de faire une première recherche de paramètres optimaux en ne prenant en compte que ceux-ci. Nous sommes conscients que nous perdons une partie de l'espace de recherche en séparant ainsi nos recherches, mais nous diminuons exponentiellement le temps de recherche.



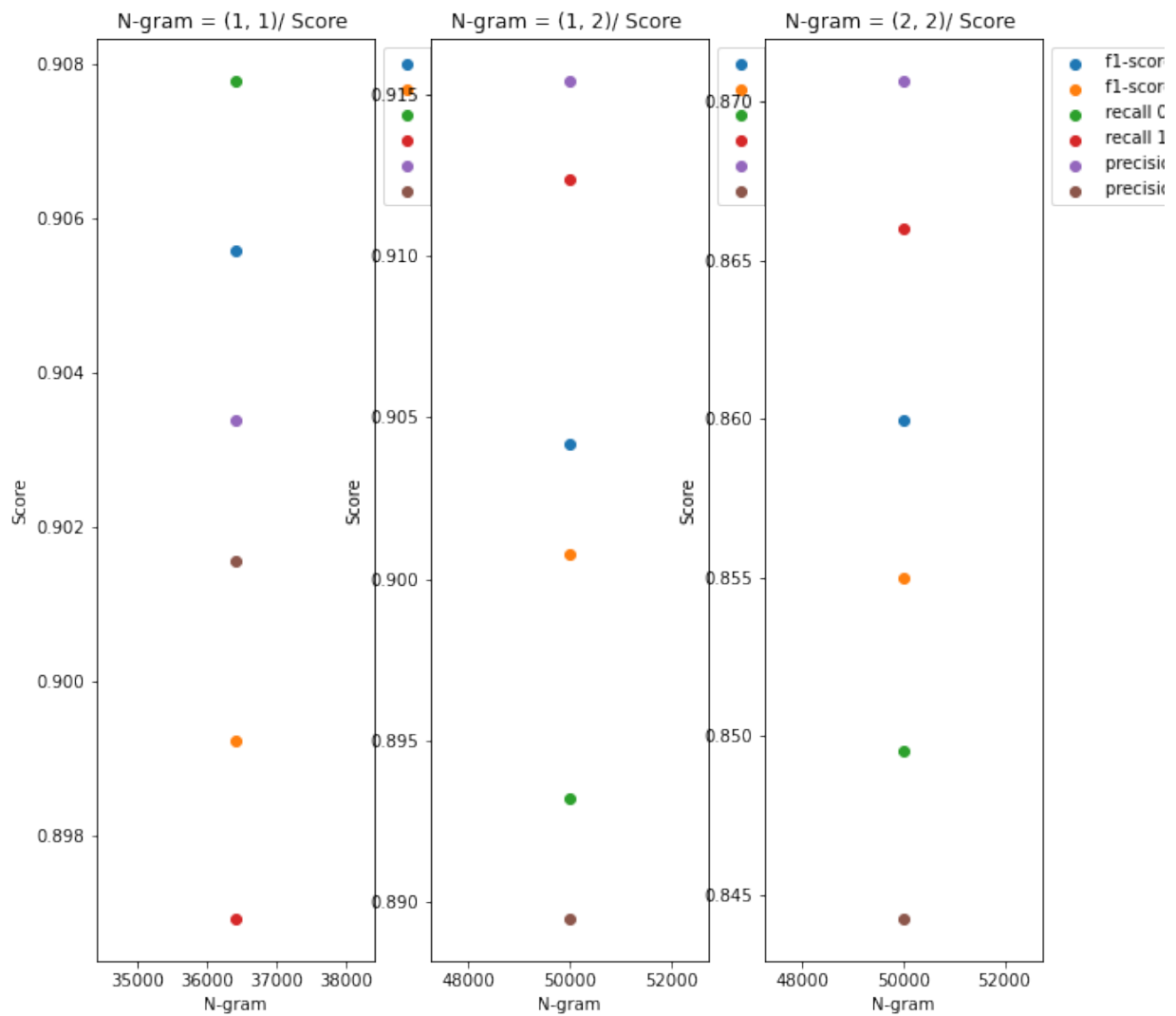


Figure 7. Variation du score suivant n-gram.

Nous utilisons le paramètre **max_features** afin d'optimiser le temps de recherche. Nous essayons de minimiser ce paramètre car il réduit la dimensionalité du problème. Il faut cependant faire attention à ne pas perdre trop de performances.

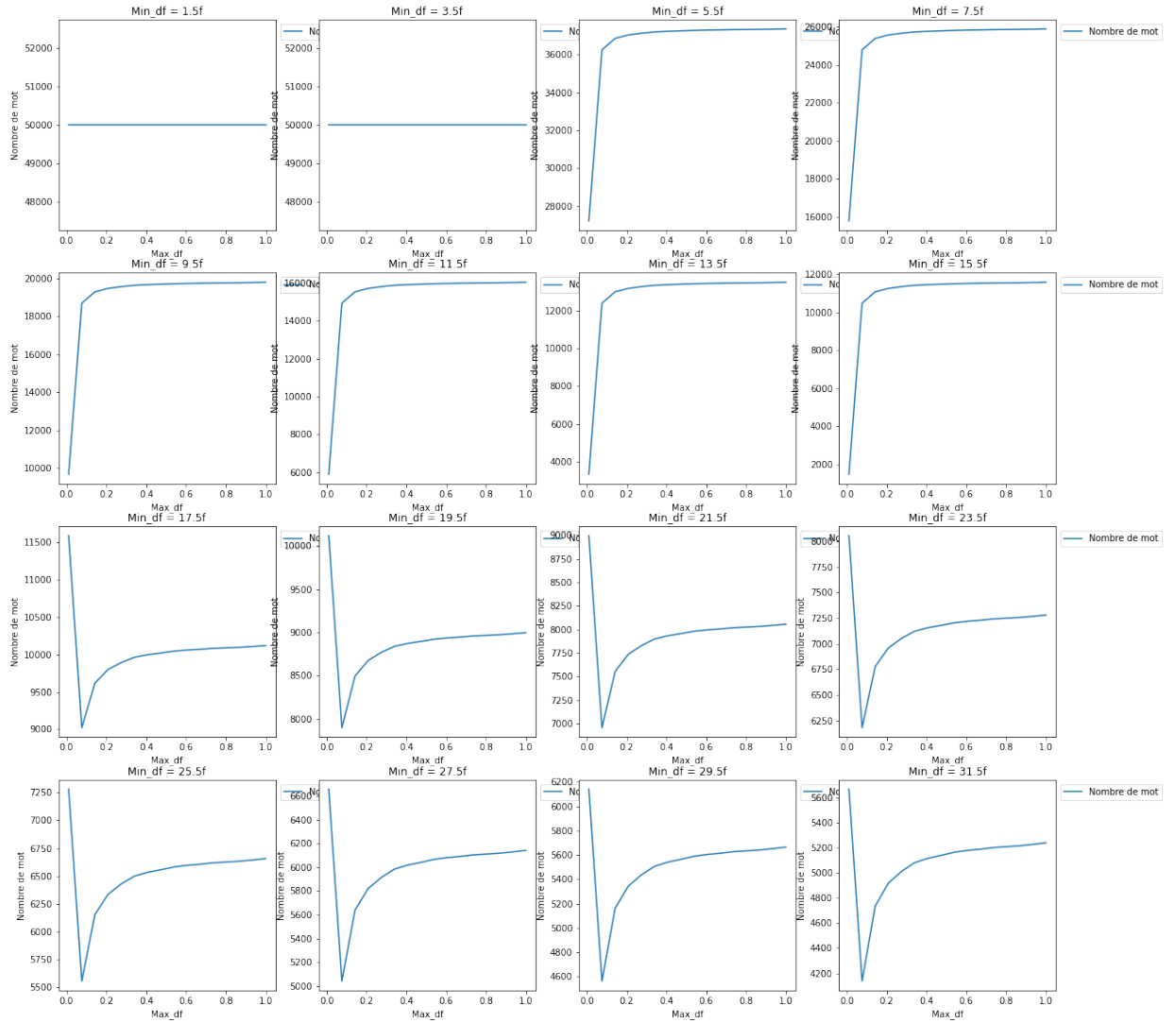


Figure 8. Variation du nombre de mots.

1.2.4. PARAMÈTRES BOOLÉENS OU DISCRETS

Nous voulons tester tous les paramètres existants de nos vectorizer. Certains comme **smooth_idf**, **sublinear_idf** influent très peu sur nos performances. D'autres comme **use_idf** sont assez importantes car il y a de grandes différences lorsque l'on utilise la normalisation par idf. Notre paramètre **use_idf** sera tout le temps activé car nous avons déjà fait notre choix lorsque nous avons choisi d'utiliser le **tfidfvectorizer**.

Les paramètres comme **ngram_range**, **stopwords**, **lowercase** et **strip_accent** sont eux plutôt influents sur nos performances. Le paramètre **ngram_range** semblent être optimal lorsque nous utilisons à la fois les unigrammes et bigrammes. En effet, les bigrammes ajoutent du contexte aux mots, ce qui est très intéressant, mais ils augmentent aussi fortement la taille du dictionnaire. Le bon compromis semble donc être l'utilisation d'unigramme et de bigramme. Pour notre problème, la suppression des **stopwords** fait baisser nos scores, cette liste de mots étant générale, on perd probablement des **stopwords** utiles à la discrimination des locuteurs (Francois Mitterrand utilise peut être beaucoup le mot "je" par exemple).

image courbe ngram lowercase etc...

1.3. Recherche des paramètres optimaux pour les modèles de machine learning

Après avoir trouvé les paramètres optimaux pour le `tfidfvectorizer` nous avons utilisé la matrice obtenue afin d'entraîner différents modèles. Nous avons le choix d'optimisation pour les paramètres suivants :

- LinearSVC : **C**, **class_weight**, **max_iter**
- MultinomialNB : **alpha**
- LogisticRegression : **C**, **class_weight**, **max_iter**

Le paramètre **class_weight** sera toujours fixé à "balanced" car la répartition des étiquettes n'est pas du tout équilibrée comme vu précédemment. Le paramètre **max_iter** sera cette fois si plus élevé, on peut se le permettre car il y a beaucoup moins de paramètres à tester pour les modèles de machine learning. L'optimisation se fait donc sur le paramètre **C** pour le SVC et la LR ou sur le paramètre **alpha** pour le MultinomialNB.

image courbe parametre C et alpha

1.4. Choix du modèle et vote

Nous avons maintenant 3 modèles optimisés ainsi qu'une bonne représentation des données. Après avoir envoyé quelques prédictions et reçu plusieurs résultats, notre modèle SVC est celui qui donne les meilleurs scores, suivis du LR. Les résultats du SVC et LR étant plutôt proche, nous essayons de faire voter ces classifieurs afin d'essayer d'améliorer notre score. Nous utilisons le module `votingclassifier` de `sklearn` et le paramètre **voting** afin de préciser que nous voulons utiliser le soft voting. Nous utilisons aussi un hard voting sur les 3 classifieurs mais il ne donne pas de bon résultat sur le jeu d'entraînement, car le f1 score du MultinomialNB est assez loin des scores du SVC et LR. Après envoi des résultats, le modèle SVC reste en tête devant le modèle de soft voting. Notre intuition sur ces résultats est que nous avons entraîné notre voting classifier sur le `tfidfvectorizer` optimisé pour le modèle SVC seulement.

1.5. Post-processing des données

Jusqu'ici, nous obtenons un f1 score sur la classe minoritaire qui est d'environ 55%. Nous avons cependant un a priori sur ce à quoi devrait ressembler notre vecteur. En effet, le texte en entrée doit être un discours entre deux présidents, il doit donc y avoir une structure par bloc. On peut donc faire un post processing des prédictions

les lissant. Pour cela, on fait passer deux algorithmes sur nos meilleures prédictions. Le premier passe sur chaque prédiction et modifie la classe de celle-ci si ses voisins font majoritairement partie de l'autre classe. On a ici 2 paramètres qui sont la **window_size** et le **threshold**. Nous apprenons ces deux paramètres grâce à une comparaison du f1 score sur le jeu de donnée. Nous obtenons en général un **threshold** de 0.6 et une **window_size** de 12

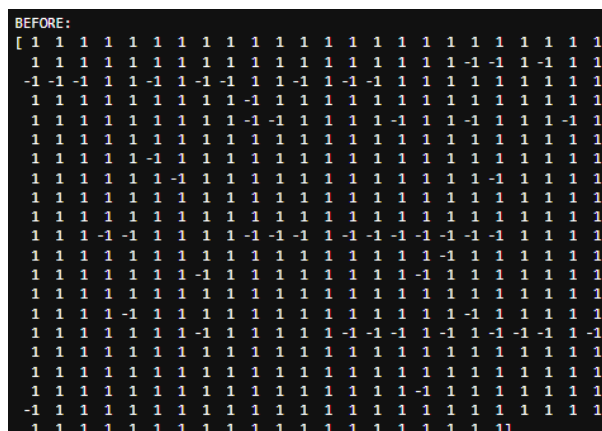


Figure 9. Avant post-processing.

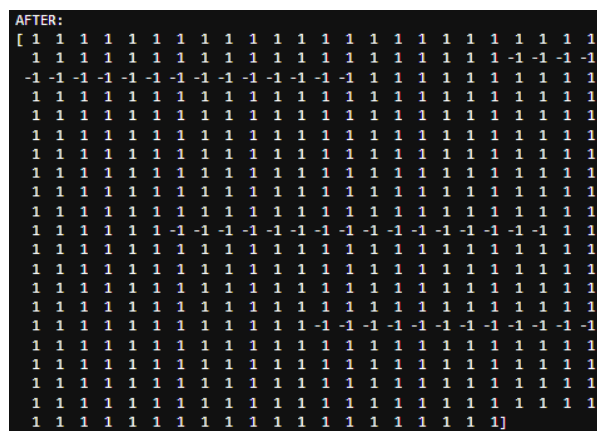


Figure 10. Après post-processing.

1.6. Conclusion

<u>TfidfVectorizer</u>	
<u>lowercase</u>	False
<u>stop_words</u>	None
<u>max_features</u>	100 000
<u>min_df</u>	1
<u>max_df</u>	0.5
<u>ngram_range</u>	(1 ,2)
<u>strip accent</u>	<u>ascii</u>
<u>use_idf</u>	<u>True</u>
<u>Smooth_idf</u>	<u>True</u>
<u>Sublinear_tf</u>	False

Figure 11. Paramètres pré-processing.

SVC	
<u>Max_iter</u>	10 000
C	100
<u>class_weight</u>	<u>balanced</u>

Pré-process	
<u>Threshold</u>	0.6
Windows	12

Figure 12. Paramètres du modèle et post-processing.

2. Jeu de données des films

Dans cette partie, on s'intéressera au jeu de données sur les avis de films. Nous avons un dossier d'apprentissage contenant deux dossiers, l'un contient des avis positifs et l'autre contient des avis négatifs. Comme pour le jeu de données précédent, le but est de classer des avis d'utilisateur contenu dans un fichier de test.

2.1. Analyse des données

Encore une fois, la première étape est de visualiser la répartition des classes dans nos données. Cette fois-ci par opposition au premier jeu de donnée, nous avons 50% de phrases à connotation positive et 50% de phrases à connotation négative. On remarque que pour ce jeu de donnée les données sont équilibrés.

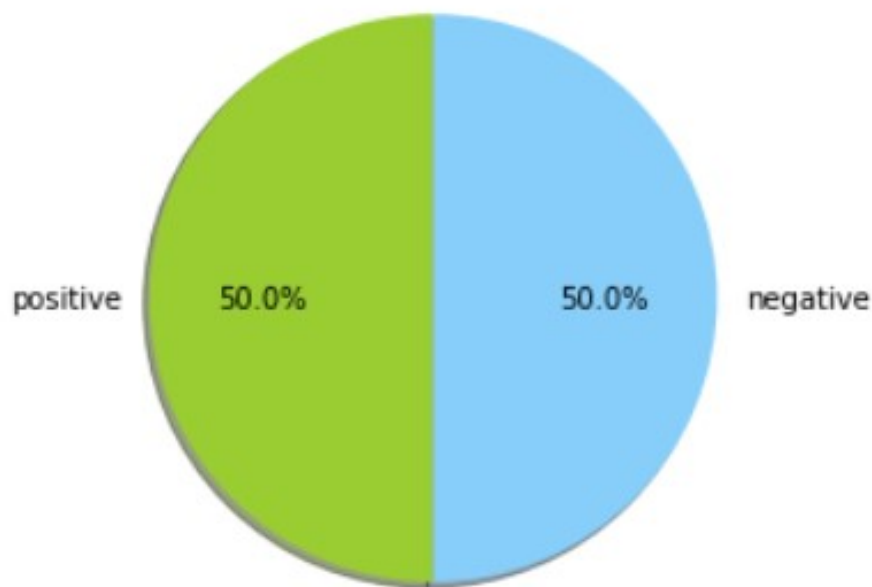


Figure 13. Répartition classes Films.

De manière analogue au jeu de donnée des présidents nous avons réalisé les mêmes campagnes d'expériences. Nous avons d'abord utilisé la librairie wordcloud afin de visualiser les mots apparaissant le plus dans notre jeu de données.

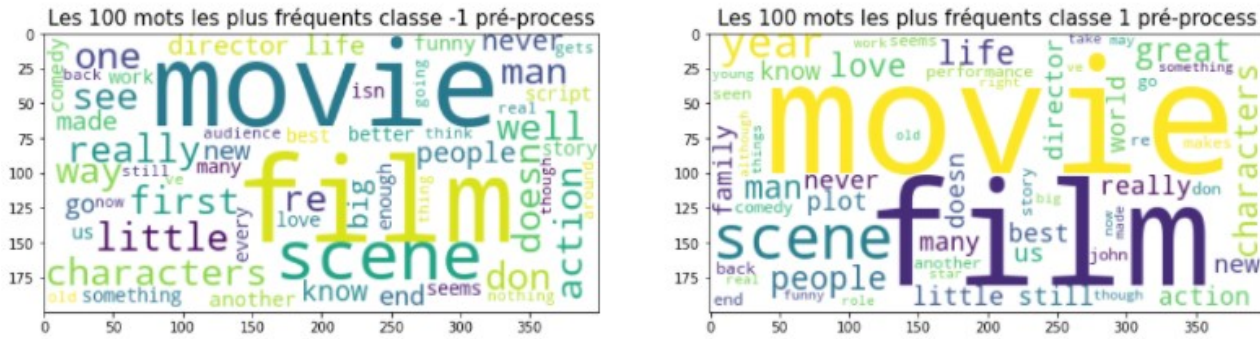


Figure 14. Top 100 mots avant pré-processing

Puis comme nous avons obtenu presque que des **stopwords**, Alors on a utilisé des listes de **stopwords** prédéfinis, et on a utilisé le paramètre **max_df** de notre vectorizer pour filtrer les mots qui n'étaient pas pertinents pour la représentation.

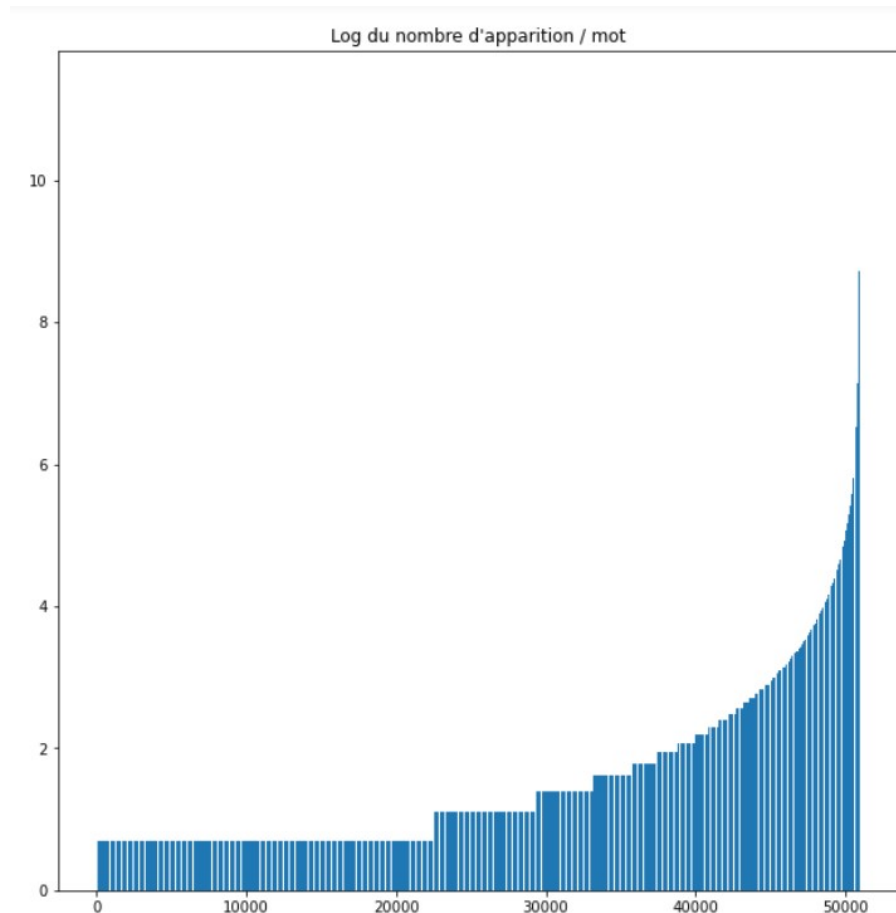


Figure 15. Loi de zipf Films.

Enfin pour confirmer nos résultats on a décidé encore une fois d'afficher les oddsratio qui nous a permit d'avoir

de meilleurs résultats.

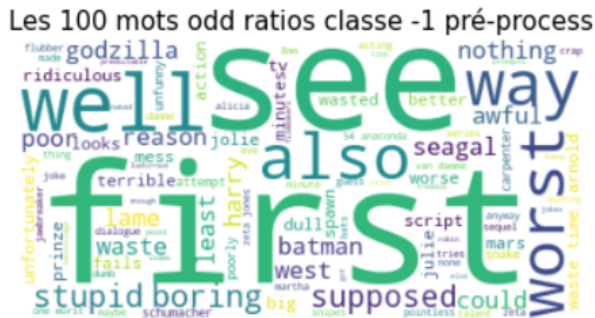


Figure 16. Top 100 mots après pré-processing négatif.



Figure 17. Top 100 mots après pré-processing positif.

2.2. Recherche des paramètres optimaux pour la représentation des données

2.2.1. REPRÉSENTATION DES DONNÉES

De nouveau, pour représenter les données la représentation qui nous a été la plus performante a été le tfidfvectorizer. La comparaison a été faite entre le countvectorizer et le tfidfvectorizer.

2.2.2. ORGANISATION

Nous avons mené une campagne d'expérience afin de trouver les paramètres optimaux pour notre problème. En effet, il faut réaliser encore une fois des fonctions de gridsearch avec différentes pipelines pour trouver les paramètres optimaux. Les métriques et les paramètres change d'un jeu de données à l'autre, c'est pourquoi on a du réentraîner nos modèles. Cette fois-ci, puisque les données étaient équilibrées l'accuracy fut suffisante pour évaluer nos données. Pour cette expérience, nous avons utilisé gridsearch sur les vectorizer (countvectorizer et tfidfvectorizer) ainsi que sur les modèles (LR, MultinomialNB et SVC).

Pour trouver les paramètres optimaux on a dû réaliser de nombreuses combinaisons. On a séparé l'ensemble des paramètres afin de gagner en temps de calcul. En effet en séparant les différents paramètres cela nous a permis d'éviter de nombreuses combinaisons qui n'étaient pas pertinentes. En ce qui concerne l'organisation, nous avons séparé nos recherches en 3 parties de la même manière que dans le jeu de données des présidents.

2.3. Recherche des paramètres optimaux pour les modèles de machine learning

Après avoir trouvé les paramètres optimaux pour le tfidfvectorizer nous avons utilisé la matrice obtenue afin d'entraîner différents modèles. Nous avons le choix d'optimisation pour les paramètres suivants :

- LinearSVC : **C**, **max_iter**
- MultinomialNB : **alpha**
- LogisticRegression : **C**, **max_iter**

L'optimisation se fait sur le paramètre **C** pour le SVC et la LR ou sur le paramètre **alpha** pour le MultinomialNB.

image courbe parametre C et alpha

2.4. Choix du modèle et vote

Nous avons maintenant 3 modèles optimisés ainsi qu'une bonne représentation des données. Après avoir envoyé quelques prédictions et reçu plusieurs résultats. Cette fois-ci, notre modèle LR est celui qui donne les meilleurs scores, suivis du SVC. Les résultats du LR et SVC étant plutôt proche, nous essayons de faire voter ces classifieurs afin d'essayer d'améliorer notre score. Mais encore une fois ces résultats n'étaient pas concluants, il faudrait optimiser séparément les différents modèles puis agréger les résultats. On retiendra que notre modèle le plus performant est LR.

2.5. Conclusion

<u>TfidfVectorizer</u>	
<u>lowercase</u>	False
<u>stop_words</u>	None
<u>max_features</u>	10 000
<u>min_df</u>	5
<u>max_df</u>	0.5
<u>ngram_range</u>	(1 ,2)
<u>strip Accent</u>	None
<u>use_idf</u>	<u>True</u>
<u>umooth_idf</u>	<u>True</u>
<u>sublinear_tf</u>	False

Figure 18. Paramètres du vectorizer.

LR	
<u>Max_iter</u>	10 000
C	100
<u>Class_weight</u>	None

Figure 19. Paramètres du modèle.