



Classification de discours présidentiel.

SIHAMDI Mostefa, BOUSBA Abdellah

UE RITAL 2020-2021, Encadrant: Vincent Guigue

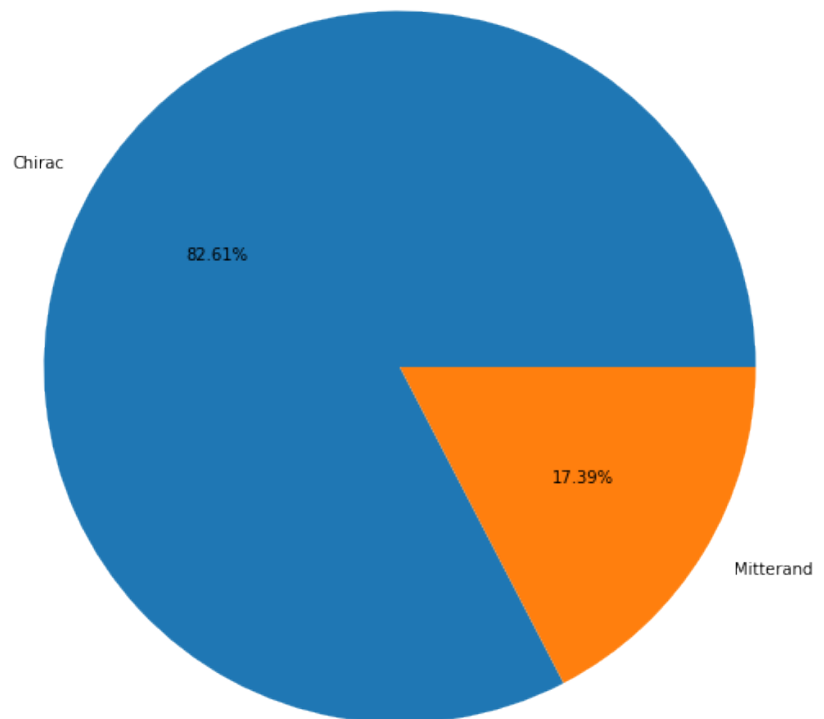
M1 DAC

I. Introduction:

Ce rapport représente les campagnes d'expériences effectuées afin d'optimiser les performances de classification de documents qui sont des discours de Jacques Chirac et François Mitterrand. En comparant plusieurs méthodes de traitement du texte et algorithmes de classification.

II. Pré-traitement :

Avant d'effectuer des changements sur le texte lui-même on remarque un déséquilibre entre les deux classes dans notre dataset train 82.61% Chirac et 17.39% Mitterrand ce qui ne nous permet pas de bien apprendre la classe minoritaire la solution la plus évidente est de faire un équilibrage par sous-échantillonnage sur la classe minoritaire, donc nos expériences seront sur notre base originale et une base équilibrée.



- Pour avoir une bonne estimation des performances de nos algorithmes nous avons choisi la métrique f1_score pour chaque classe (plus précisément celle de Mitterrand).

Equilibrage des données :

	Chirac Mean	Chirac Std	Mitterrand Mean	Mitterrand Std
SVM	0.934872	0.001887	0.487983	0.017170
Naive Bayes	0.931559	0.002962	0.518929	0.010911
Reg Logistique	0.941576	0.002016	0.488863	0.017523

Score F1 par cross validation sans équilibrage CountVectorizer

	Chirac Mean	Chirac Std	Mitterrand Mean	Mitterrand Std
SVM	0.742652	0.011267	0.738125	0.009267
Naive Bayes	0.769704	0.009202	0.776069	0.003987
Reg Logistique	0.768454	0.008353	0.760677	0.006033

Score F1 par cross validation avec équilibrage CountVectorizer

Combinaison optimale pour le traitement du texte :

	to_lower	keep_punct	stemming	keep_stopwords	f1_C	f1_M
0	1.0	1.0	1.0	1.0	0.748263	0.701984
1	1.0	1.0	1.0	0.0	0.738003	0.689421
2	1.0	1.0	0.0	1.0	0.752391	0.704109
3	1.0	1.0	0.0	0.0	0.736992	0.687667
4	1.0	0.0	1.0	1.0	0.747347	0.700097
5	1.0	0.0	1.0	0.0	0.733292	0.684401
6	1.0	0.0	0.0	1.0	0.752391	0.704109
7	1.0	0.0	0.0	0.0	0.736304	0.688104
8	0.0	1.0	1.0	1.0	0.748263	0.701984
9	0.0	1.0	1.0	0.0	0.740167	0.690609
10	0.0	1.0	0.0	1.0	0.765724	0.720292
11	0.0	1.0	0.0	0.0	0.761340	0.717465
12	0.0	0.0	1.0	1.0	0.747347	0.700097
13	0.0	0.0	1.0	0.0	0.735036	0.686184
14	0.0	0.0	0.0	1.0	0.765724	0.720292
15	0.0	0.0	0.0	0.0	0.755064	0.710939

GridSearch SVM CountVectorizer

	to_lower	keep_punct	stemming	keep_stopwords	f1_C	f1_M
0	1.0	1.0	1.0	1.0	0.773267	0.742249
1	1.0	1.0	1.0	0.0	0.762244	0.732706
2	1.0	1.0	0.0	1.0	0.779882	0.755379
3	1.0	1.0	0.0	0.0	0.768583	0.743164
4	1.0	0.0	1.0	1.0	0.774005	0.743317
5	1.0	0.0	1.0	0.0	0.761406	0.731765
6	1.0	0.0	0.0	1.0	0.779882	0.755379
7	1.0	0.0	0.0	0.0	0.765246	0.734906
8	0.0	1.0	1.0	1.0	0.773267	0.742249
9	0.0	1.0	1.0	0.0	0.765602	0.734954
10	0.0	1.0	0.0	1.0	0.791745	0.763886
11	0.0	1.0	0.0	0.0	0.788132	0.760962
12	0.0	0.0	1.0	1.0	0.774005	0.743317
13	0.0	0.0	1.0	0.0	0.771464	0.737993
14	0.0	0.0	0.0	1.0	0.791745	0.763886
15	0.0	0.0	0.0	0.0	0.786906	0.758507

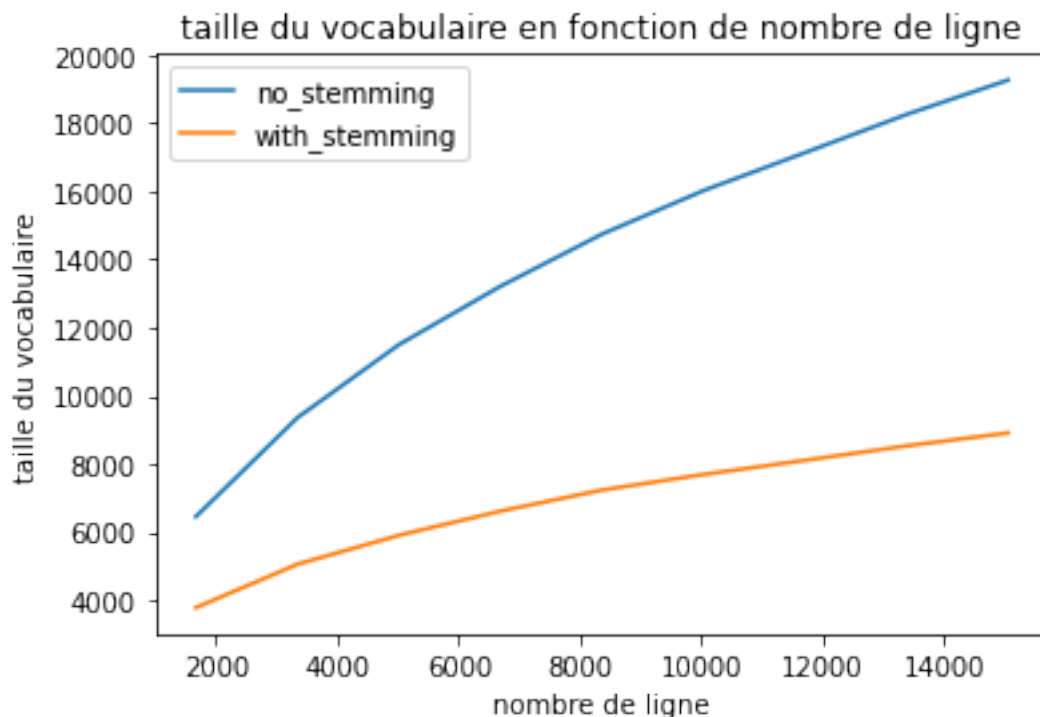
GridSearch Naïve Bayes CountVectorizer

	to_lower	keep_punct	stemming	keep_stopwords	f1_C	f1_M
0	1.0	1.0	1.0	1.0	0.773720	0.731282
1	1.0	1.0	1.0	0.0	0.759071	0.713246
2	1.0	1.0	0.0	1.0	0.779173	0.733627
3	1.0	1.0	0.0	0.0	0.766850	0.721885
4	1.0	0.0	1.0	1.0	0.772801	0.728863
5	1.0	0.0	1.0	0.0	0.754840	0.708081
6	1.0	0.0	0.0	1.0	0.779173	0.733627
7	1.0	0.0	0.0	0.0	0.763067	0.716614
8	0.0	1.0	1.0	1.0	0.773720	0.731282
9	0.0	1.0	1.0	0.0	0.763770	0.719341
10	0.0	1.0	0.0	1.0	0.790224	0.749879
11	0.0	1.0	0.0	0.0	0.782466	0.743849
12	0.0	0.0	1.0	1.0	0.772801	0.728863
13	0.0	0.0	1.0	0.0	0.765507	0.719533
14	0.0	0.0	0.0	1.0	0.790224	0.749879
15	0.0	0.0	0.0	0.0	0.774856	0.736538

GridSearch RegLog CountVectorizer

On peut clairement voir que la suppression des stopwords affecte le score négativement ce qui peut être expliqué par le fait que le concept de stopwords peut varier d'une base à une autre, dans notre cas la suppression supprime des informations utiles pour la classification, de même pour les minuscules, par contre même si le stemming ne contribue pas trop pour le score on testera ensuite son effet sur la taille du vocabulaire pour décider ensuite si ça vaut le coût de le garder ou pas.

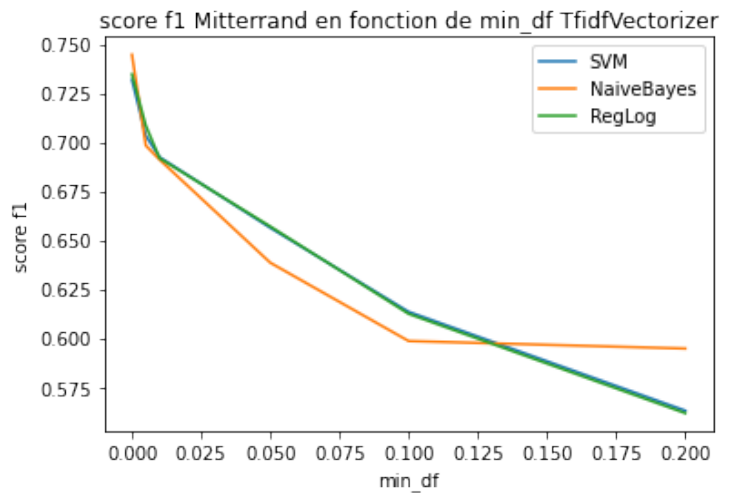
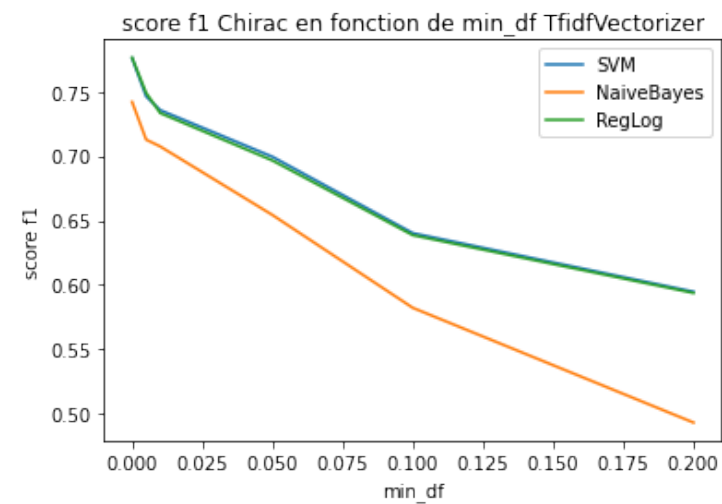
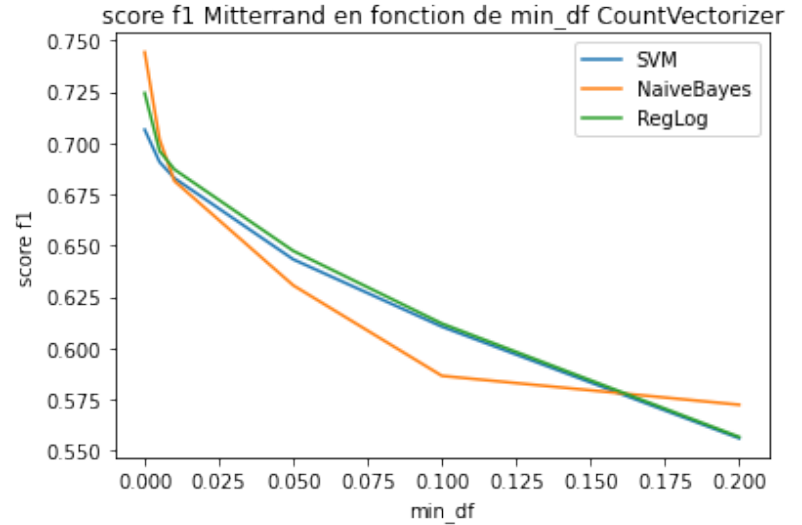
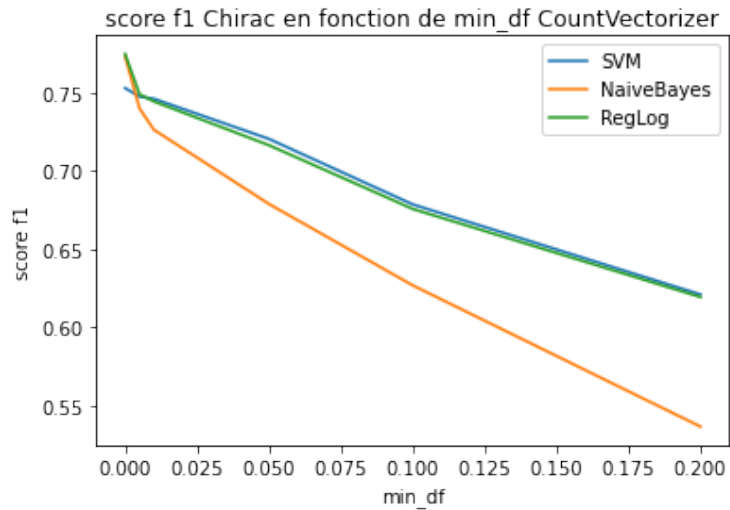
Effet du stemming :



On remarque que le stemming diminue significativement la taille du vocabulaire ce qui nous permettra d'avoir un meilleur temps d'exécution.

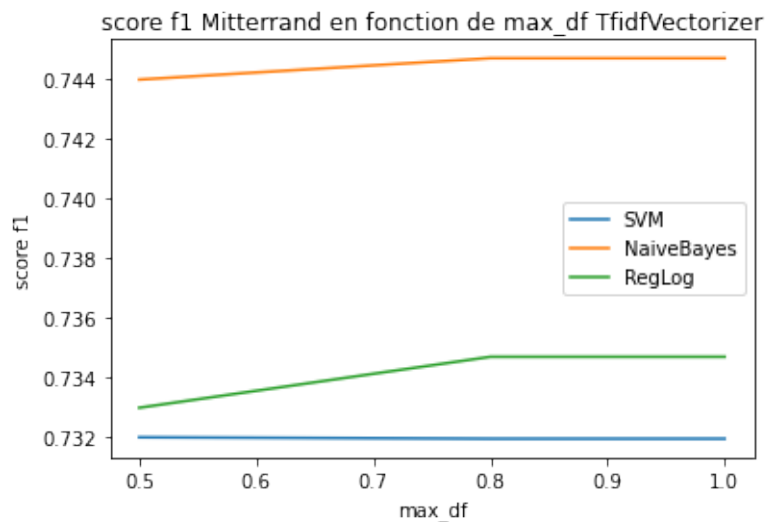
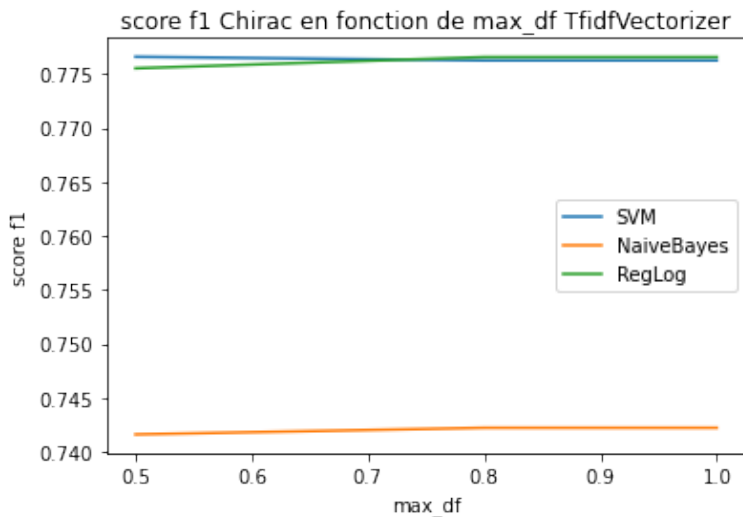
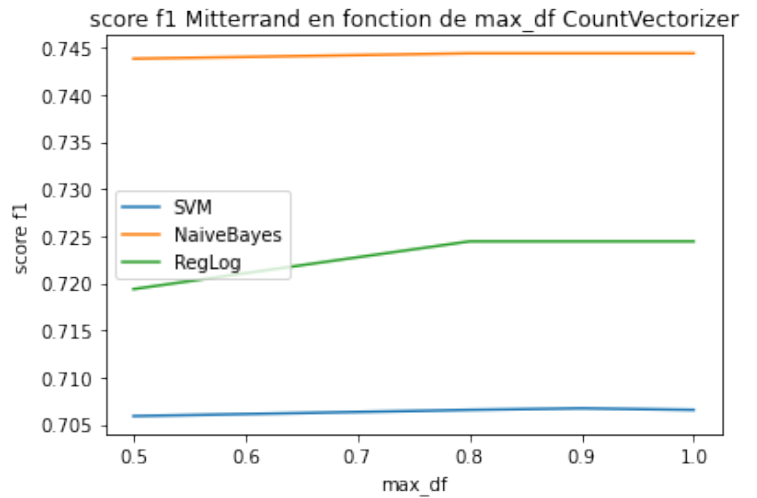
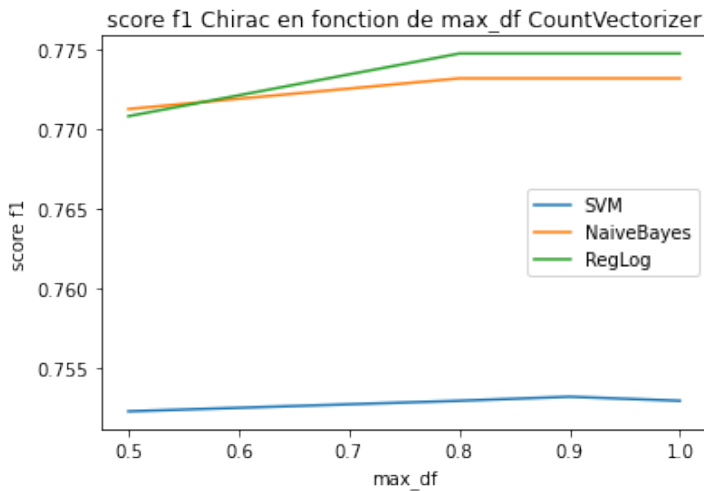
Fréquence documentaire des mots :

min_df :



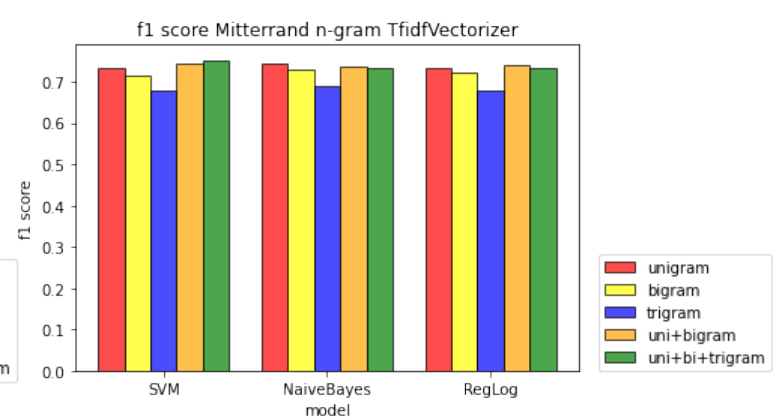
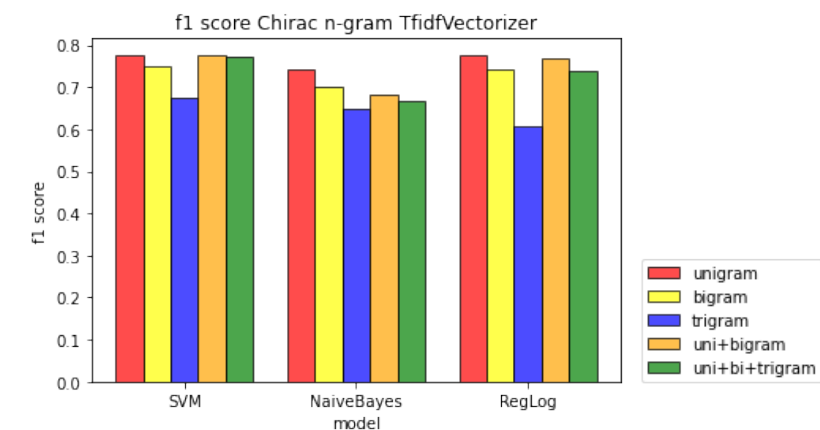
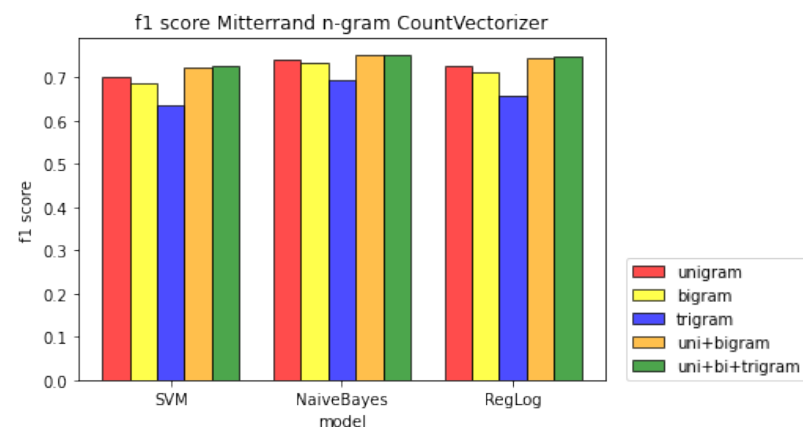
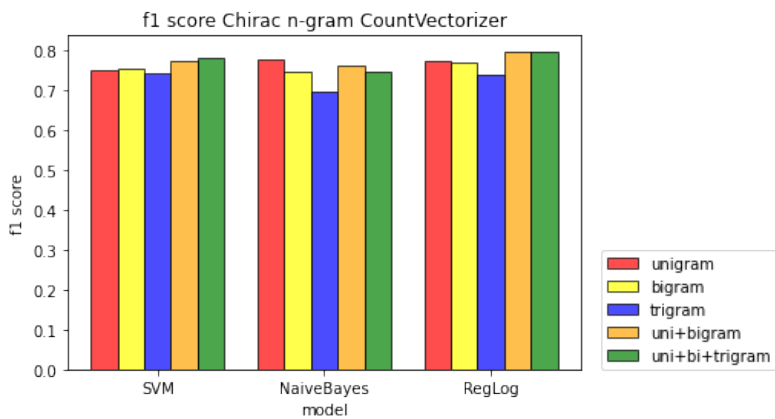
Après avoir tester la suppression des mots en fonction de la fréquence documentaire minimum, on remarque des résultats similaires pour CountVectorizer et TfidfVectorizer, en plus ce traitement n'améliore pas le score.

max_df :



Après avoir tester la suppression des mots en fonction de la fréquence documentaire maximum, on remarque des résultats similaires pour CountVectorizer et TfidfVectorizer sauf qu'on préfère utiliser CountVectorizer au lieu de TFIDF avec NaiveBayes, en plus pour le classe Mitterrand on remarque que NaiveBayes a des résultats plus intéressante, la suppression des mots avec une grande fréquence documentaire ne contribue pas pour le score .

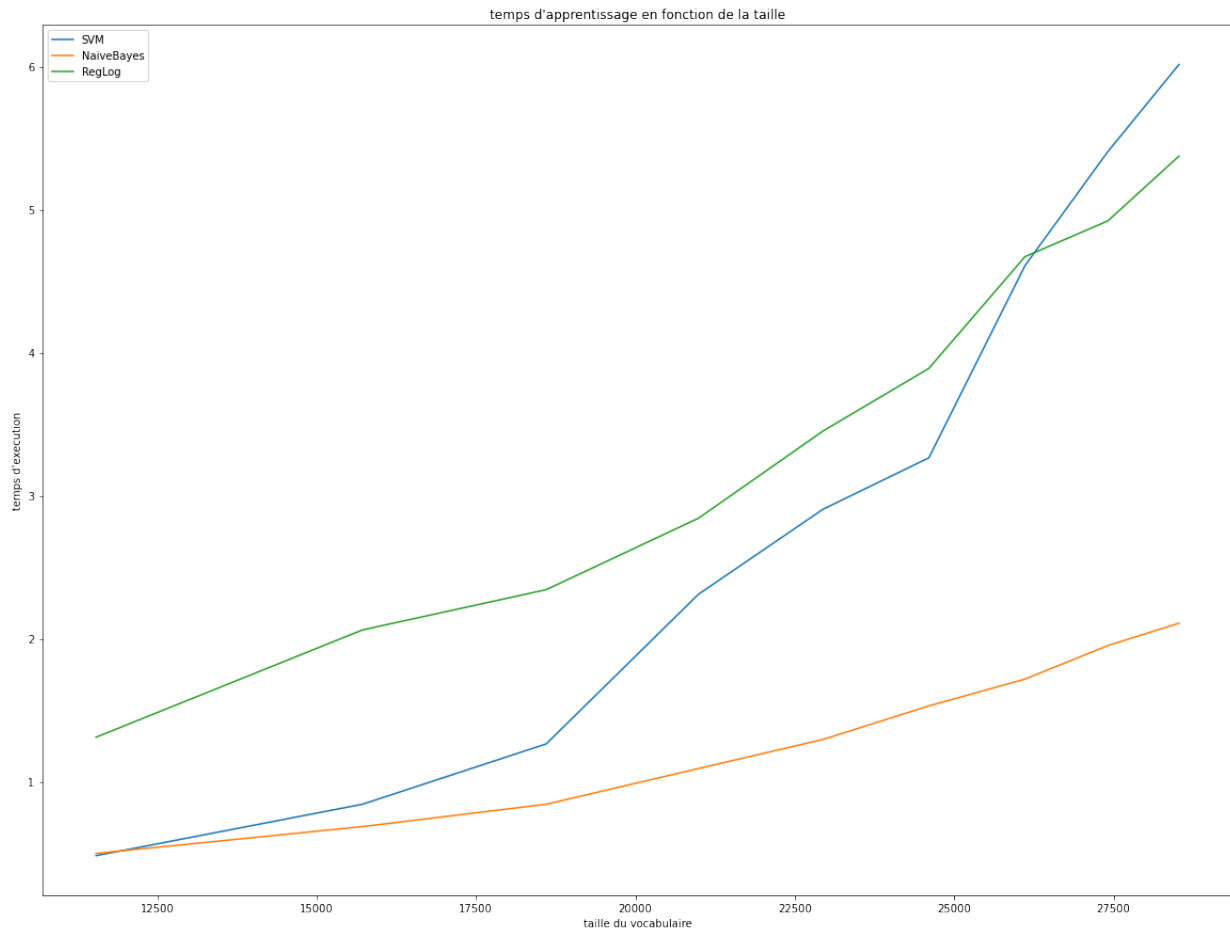
n-gram :



Concernant les tests au-dessus on ne remarque que chaque algorithme est affecté différemment, pour le **SVM** la meilleure combinaison est le mélange entre uni bi et trigram en plus le TfidfVectorizer retourne des résultats plus élevés, pour le **Naive Bayes** si on utilise CountVectorizer un mélange entre uni et bigram paraît mieux, mais si on utilise TfidfVectorizer il suffit de prendre des unigram sachant que CountVectorizer est mieux approprié pour cet algorithme, et les mêmes remarques peuvent être dites sur la **Régression Logistique**

III. Classifieur:

Temps d'exécution:

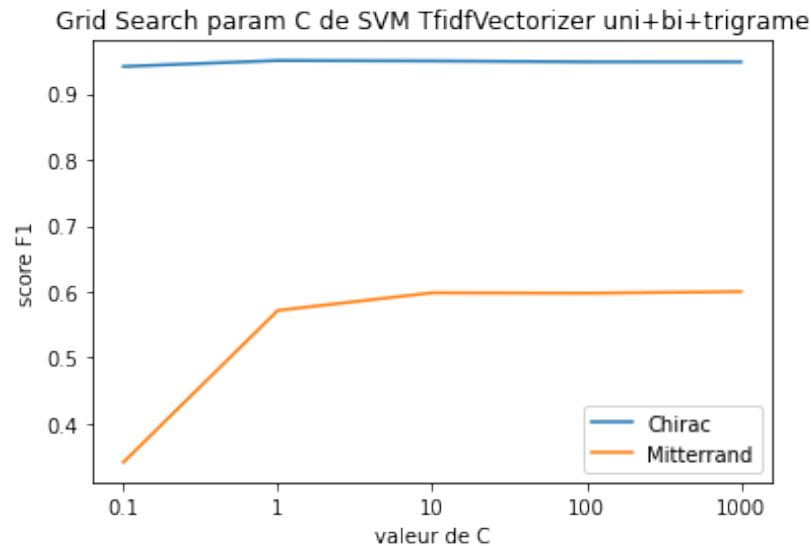


Après avoir essayer plusieurs combinaisons on arrive à la conclusion que le modèle avec meilleurs résultats est :

- **SVM : TfidfVectorizer + unigram+bigram+trigram**

Même si le modèle Naive Bayes est plus rapide, on a décidé d' utiliser le SVM dans le but de gagner de la précision au lieu du temps. Afin d'avoir une meilleure estimation de performance nous avons découper notre dataset train en 70% train et 30% test en regardant le id des documents au lieu d'utiliser une cross validation ou split train test avec shuffle pour garder la même structure des données. Et après exécution voici les résultats obtenus :

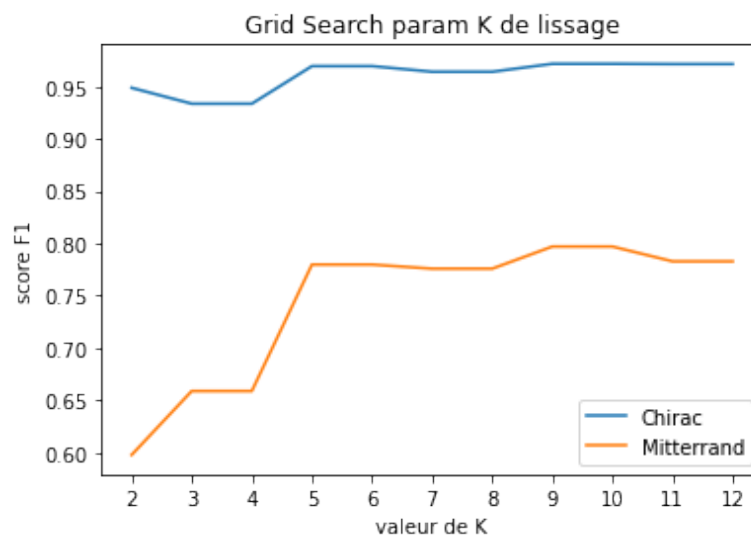
Après avoir fait un équilibrage de train notre score f1_score Mitterrand : 0.478, Chirac : 0.852, par contre en laissant les données déséquilibrées on obtient un meilleur score f1_score **Mitterrand : 0.597**, **Chirac : 0.948**, on peut expliquer ça par la taille des données train qui est très petite dans le cas d'équilibrage. Pour aller plus loin on peut s'en servir de nos connaissances sur le jeu de données pour faire un lissage sur les étiquettes.



Le changement du paramètre C n'affecte pas le score de Chirac mais pour C=100 on améliore bien le score.

IV. Post- traitement:

Nous observons que les données sont sous formes de blocs. C'est-à-dire si un président parle, il dit plusieurs lignes à la fois, sachant que Chirac parle plus que Mitterrand, donc on fait une transformation basant sur les K plus proche voisin avec un poids de 2.5 pour le nombre des étiquettes Mitterrand.



En faisant un GridSearch pour le meilleur param de K on remarque que 9 était la meilleure valeur avec le nouveau score après lissage : **Mitterrand : 0.796, Chirac : 0.971.**

V. Conclusion:

Après avoir tester plusieurs choses on arrive à la conclusion que travailler avec des données déséquilibrer est un vrai challenge, nous avons retenu que les stopwords ne sont pas toujours considéré comme bruit, les lettres majuscules contient aussi de l'information, le stemming n'apporte pas une grande différence mais diminue la taille du vocabulaire. De plus, le choix du classifieurs est très important après les hyperparamètre les plus optimaux et la représentation correcte des données. Enfin la connaissance du format de notre dataset est très importante comme on a vu que grâce au post-traitement le score à augmenter significativement.