

# Rapport : Projet TAL – RI

## NADLER Bastian

### I) Code python :

#### 1) Introduction

Il y a 6 fichiers python différents dans lesquels j'ai appliqué les différentes méthodes présentées par les TP. Ainsi dans le code présenté, il y a :

- projet\_tdidf.ipynb correspond aux méthodes utilisées dans le 1<sup>er</sup> Tp
- projet\_CNN\_FP.ipynb et projet\_CNN\_FM.ipynb correspondent aux méthodes utilisées dans le 2<sup>er</sup> Tp
- projet\_BILSTM\_FM.ipynb et projet\_BILSTM\_FP.ipynb correspondent aux méthodes utilisées dans le 3<sup>er</sup> Tp (Les fichiers FP et FM se distinguent uniquement par le modèle utilisé, FM utilise [1] et FP [2])
- projet\_transformer.ipynb correspond aux méthodes utilisées dans le 4<sup>er</sup> Tp

Un large éventail d'algorithmes a été utilisé pour classer automatiquement les films en fonction de leur genre avec un apprentissage sur leur synopsis et leur titre. Ces algorithmes de classification seront présentés par ordre croissant de précision sur la prédiction du genre. À noter que pour la plupart des apprentissages, différentes approches ont été envisagées et la validation croisée est toujours mise en place quand cela est envisageable *i.e.* quand cela était faisable en un temps raisonnable. Ainsi la validation croisée est notamment faite dans les 3 premiers codes ci-dessus. Dans projet\_transformer elle n'est pas réalisée car elle prend beaucoup trop de temps et de puissance de calcul. De plus, je me suis concentré sur l'utilisation de tous les algorithmes présents dans les tps afin de comparer leur efficacité, en testant peu d'architectures différentes.

#### 2) CNN

Cet algorithme est celui ayant eu les moins bons résultats en terme de précision. Pour les modèles de plongements de mots pré-entraînés, 2 modèles différents ont été utilisés trouvables sur les sites présents en bibliographie [1] et [2]. Le modèle de [1] donne des résultats moins bons que ceux de [2]. Notamment le modèle [2] donne des résultats similaires à un code n'utilisant pas de modèle (à 0.02 près). Nous obtenons les résultats suivants :

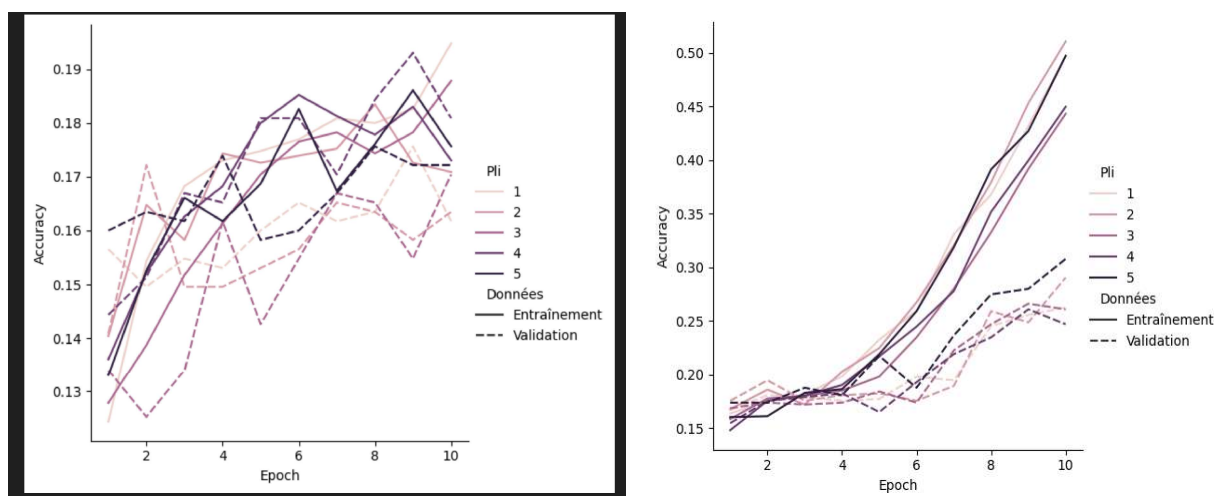
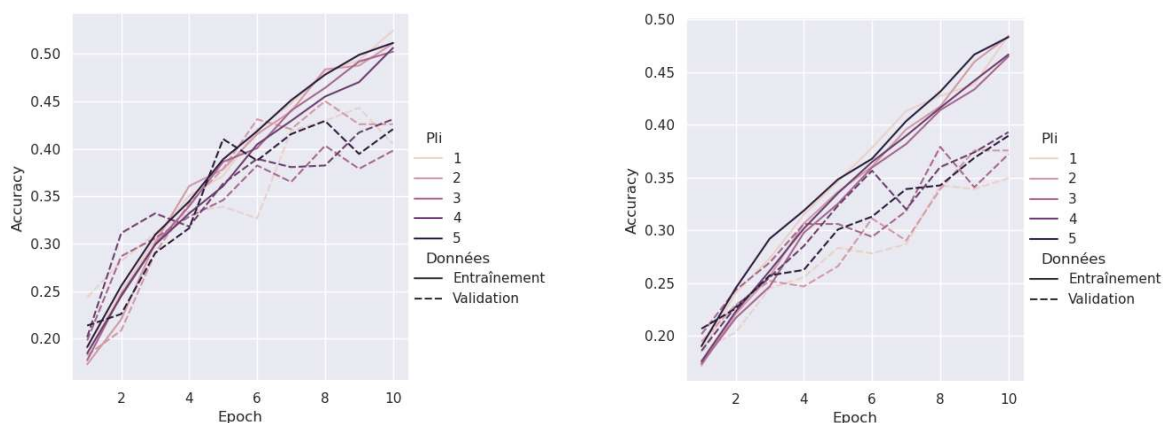


Figure 1 et 2 : Accuracy de l'algorithme CNN. À gauche avec l'utilisation du modèle [1] de plongement de mots pré-entraîné, et à droite sans.

Nous avons en moyenne 0,16 de précision avec l'utilisation du modèle [1] et 0,27 sans (ou en utilisant le modèle de [2]). Sachant que les 2 modèles de plongements de mots pré-entraînés fonctionnent pour l'algorithme BILSTM, il est probable que l'algorithme CNN soit mauvais pour résoudre notre problème. Pour améliorer les résultats, il faut très probablement un plus grand nombre de données. De plus, il faut probablement modifier son architecture *i.e.* changé le nombre de couches, la taille des filtres, les fonctions d'activation et les techniques de régularisation en fonction. A noté que j'ai tout de même essayé de changer l'architecture, mais cela n'a pas donné de bien meilleurs résultats.

### 3) **BILSTM**

De même que précédemment, les tests ont été effectués avec les 2 modèles de plongements de mots pré-entraînés présent en bibliographie. Néanmoins, les 2 modèles donnent sensiblement les mêmes résultats ici (à  $\pm 0.02$  de précision). Ainsi en utilisant LSTM et GRU, on trouve comme résultats :



**Figure 3 et 4 : Accuracy de l'algorithme BILSTM. A gauche avec LSTM et à droite avec GRU**

En moyenne, on a une précision de 0,37 pour GRU et 0,42 pour BILSTM. On voit qu'avec cet algorithme on obtient une bien meilleure précision qu'avec l'algorithme CNN. Il y a cependant plusieurs pistes pour améliorer ce résultat. Tout d'abord, comme pour CNN, un plus grand nombre de données permettrait une plus grande précision. Un prétraitement des données (tels que la désuffixation) pourrait aussi améliorer les résultats. Il faudrait aussi chercher à modifier l'architecture du réseau afin d'en trouver une plus optimale. Je n'ai pas réalisé de comparaison avec d'autre architecture car je n'ai pas réussi à construire d'architecture plus optimale.

### 4) **tdidf**

Dans cette partie, plusieurs algorithmes différents ont été comparés, avec à chaque fois des paramètres différents. Comme algorithmes utilisés, nous avons : Random forest, KNN, LR, CART, Multinomial NB, Baseline. Des tests ont été réalisés avec une utilisation de la désuffixation, l'application de différents seuils de fréquence pour les tokens à conserver dans le vocabulaire, utilisation de différentes méthodes de pondération. Néanmoins, la pondération (sublinear\_tf, use\_idf, smooth\_idf dans TfidfVectorizer) ne change pas sensiblement les résultats (l'accuracy est modifiée de l'ordre de 0.01). Ainsi nous avons les résultats suivants :

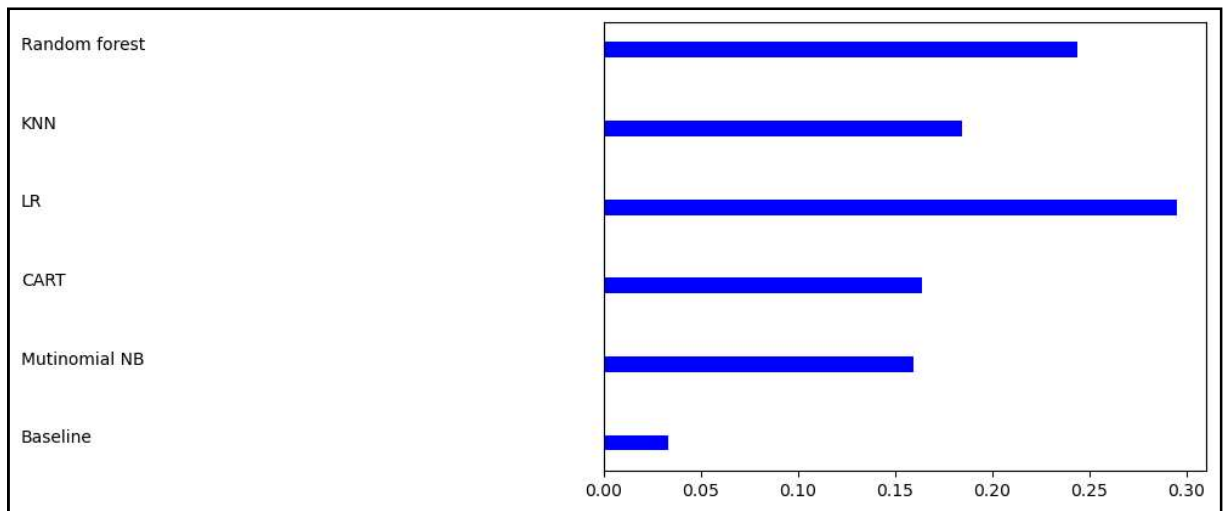


Figure 5 : Fréquence de 0,1 et sans désuffixation

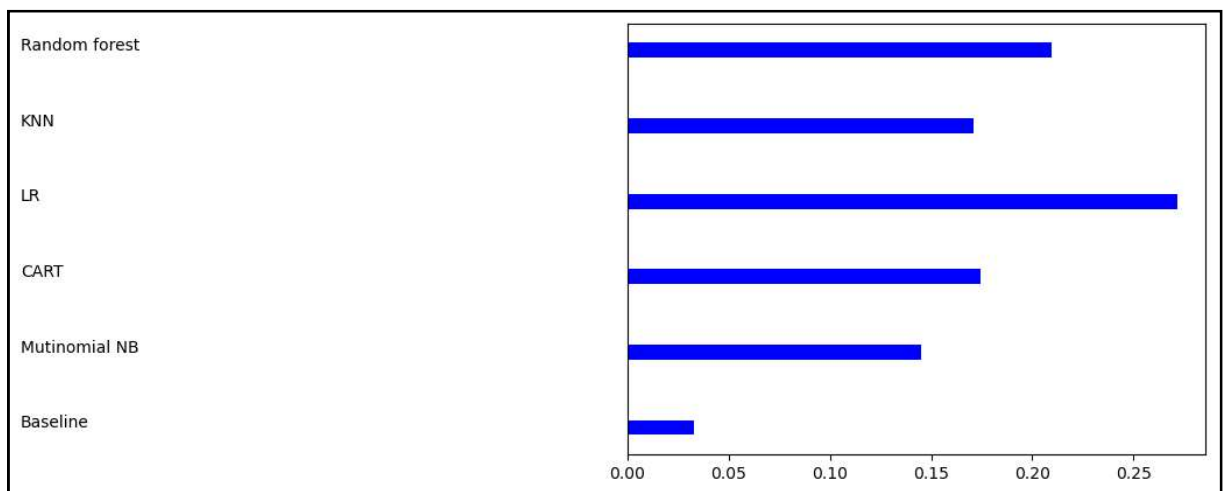


Figure 6 : Fréquence de 0,1 et avec désuffixation

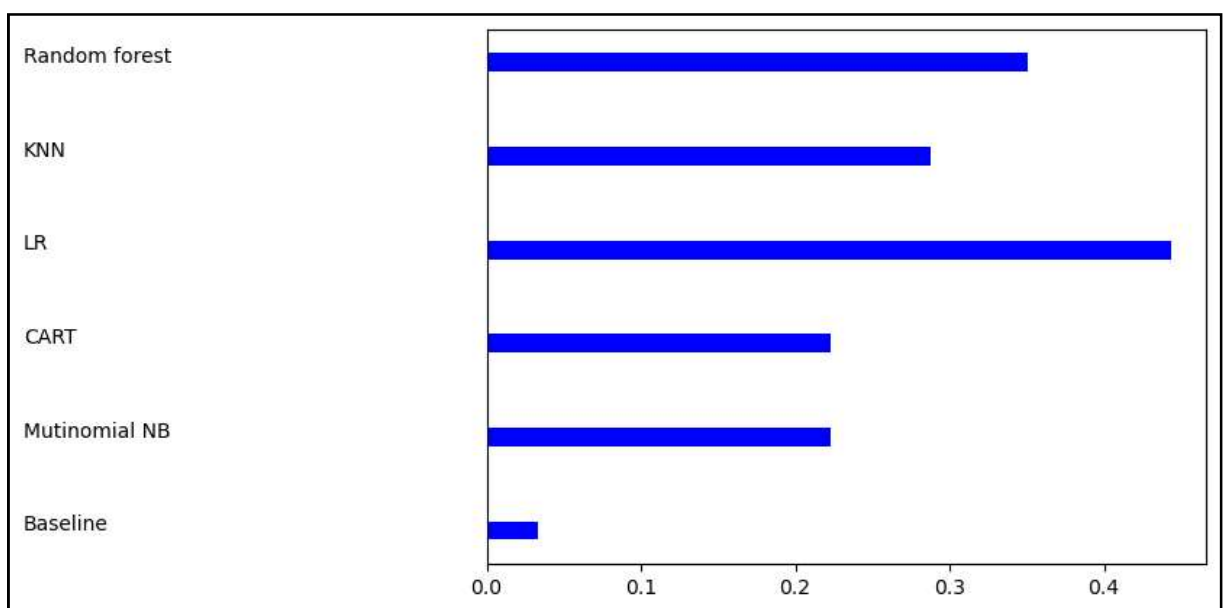


Figure 7 : Fréquence de 0,01 et sans désuffixation

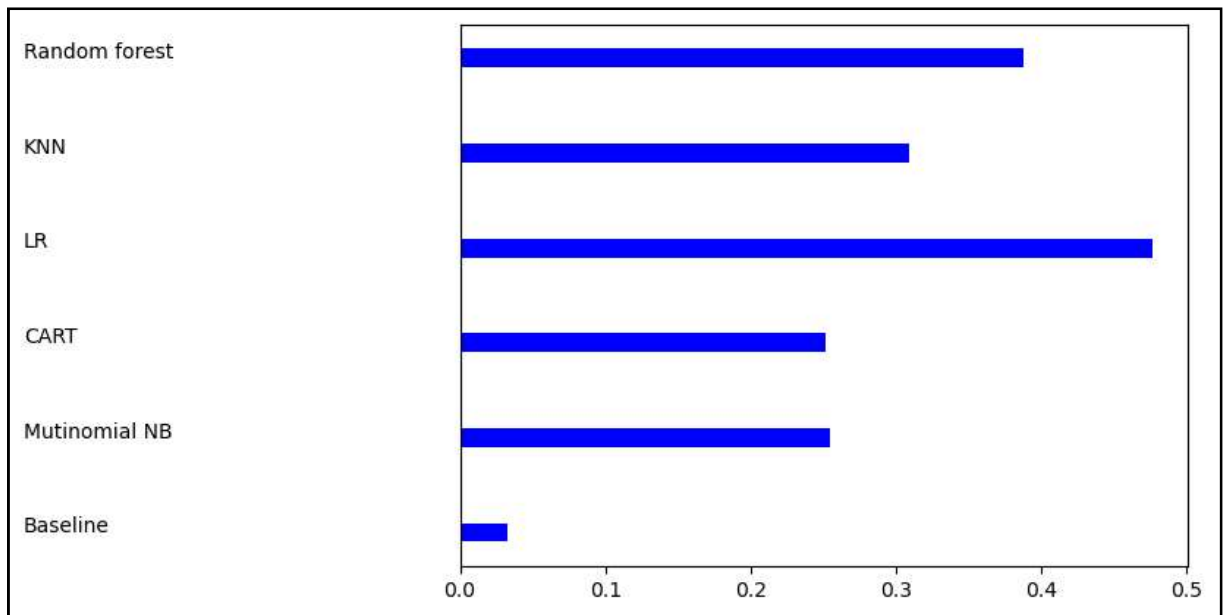


Figure 8 : Fréquence de 0,01 et avec désuffixation

Nous pouvons voir que l'algorithme le plus efficace est tout le temps LR. Il est le plus précis avec une fréquence de 0,01 et la désuffixation, atteignant une précision de 0,48 (0.43 sans désuffixation). Pour améliorer cette précision, il faudrait idéalement plus de données et veiller à un équilibrage des classes (ou gérer le déséquilibre).

## 5) Transformer

La dernière méthode utilisée est le Transformer. C'est celle qui donne les meilleurs résultats avec une accuracy de 0,60. Pour celle-ci on utilise le modèle « xlm-roberta-base » pour la tokenisation des données, qui est le plus utilisé pour la langue française. On a ainsi les résultats suivant sur le fichier allocine\_train.csv:

`preds_output.metrics`

```
{'test_loss': 1.1796503067016602,
'test_accuracy': 0.6139130434782609,
'test_runtime': 19.4273,
'test_samples_per_second': 29.598,
'test_steps_per_second': 3.706}
```

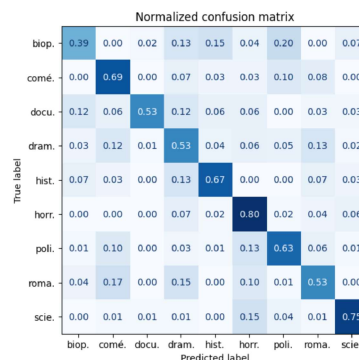


Figure 9 et 10 : Résultats du transformer lors du test et matrice de confusion

Ainsi comme c'est la méthode ayant les meilleurs résultats, le fichier allocine\_test\_enrichi.csv a été créé grâce à cette méthode. Pour améliorer la précision, il faudrait beaucoup plus de données, et modifier le learning rate, le nombre de batch, le nombre d'époque et le poids pourrait améliorer la précision. Néanmoins j'ai fait peu de test de modification de paramètre car les transformers prennent énormément de temps.

## **6) Conclusion :**

De nombreuses méthodes ont été utilisées. La plus efficace utilise les transformer est atteint une précision de 0,60. En général, les pistes d'améliorations sont l'augmentation de données, la modification de la structure des réseaux et un prétraitement des données. Il faudrait aussi annuler l'impact négatif du déséquilibre des classes sur la précision des algorithmes. Pour les algorithmes utilisant des modèles provenant du net (xlm-roberta-base, les prolongements de mots pré-entraînés) on pourrait en essayer d'autre pour voir si cela peut améliorer les résultats. Néanmoins pour les prolongements de mots pré-entraînés, les 2 utilisés lors du projet sont assez fournis (entre 1.5Gb et 2.2Gb) et il y a peu de chance qu'un autre puisse améliorer drastiquement les résultats. Et aucun autre modèle que xlm-roberta-base n'a été testé pour les transformers car cela prend trop de temps et de puissance de calculs.

## II) TAL

Le document TAL a été créé et ajusté pour correspondre à l'image donnée dans l'intitulé du projet. Le résultat final devrait être cela :

localhost:8983/solr/projet/browse?q=cheval 80 %

cheval Envoyer Reset

3 results found in 4ms Page 1 of 1

**Filtrer par ...**

**Réalisateur**

Andrew Haigh (1)  
Eric Toledano, Olivier Nakache (1)  
Nabil Ayouch (1)

**Année**

2008 (1)  
2012 (1)  
2017 (1)

**Nationalité**

Américain (1)  
Français (1)  
Marocain, français (1)

**Langues**

Anglais (1)  
Arabe (1)  
Français (1)

**Genre**

drame (2)  
comédie (1)

**genre\_predit**

drame (2)  
comédie (1)

**annee\_prod:** 2012  
**langues:** Arabe  
**nationalite:** Marocain, français, belge  
**realisateurs:** Nabil Ayouch  
**synopsis:** Yassine a 10 ans lorsque le Maroc émerge à peine des années de plomb. Sa mère, Yemma, dirige comme elle peut toute la famille. Un père dépressif, un frère à l'armée, un autre presque autiste et un troisième, Hamid, petit calé du quartier et protecteur de Yachine. Quand Hamid est emprisonné, Yachine enchaîne les petits boulots. Pour les sortir de ce marasme où règnent violence, misère et drogue, Hamid, une fois libéré et devenu islamiste radical pendant son incarcération, persuade Yachine et ses copains de rejoindre leurs "frères". L'Imam Abou Zoubair, chef spirituel, entame alors avec eux une longue préparation physique et mentale. Un jour, il leur annonce qu'ils ont été choisis pour devenir des martyrs ...  
**titre:** Les **Chevaux** de Dieu  
**genre:** drame  
**genre\_predit:** drame  
**url:** 6cb1b6e4-3f7e-48c0-adfc-406f10b926ba

**annee\_prod:** 2017  
**langues:** Anglais  
**nationalite:** Américain  
**realisateurs:** Andrew Haigh  
**synopsis:** Charley Thompson a quinze ans et a appris à vivre seul avec un père inconstant. Tout juste arrivé dans l'Oregon, le garçon se trouve un petit boulot chez un entraîneur de **chevaux** et se prend d'affection pour Lean on Pete, un pur-sang en fin de carrière. Le jour où Charley se retrouve totalement livré à lui-même, il décide de s'enfuir avec Lean on Pete, à la recherche de sa tante dont il n'a qu'un lointain souvenir. Dans l'espoir de trouver enfin un foyer, ils entament ensemble un long voyage ...  
**titre:** La Route sauvage (Lean on Pete)  
**genre:** drame  
**genre\_predit:** drame  
**url:** 622cb4da-e50e-48f4-8258-27a778dd443a

**annee\_prod:** 2008  
**langues:** Français  
**nationalite:** Français  
**realisateurs:** Eric Toledano, Olivier Nakache  
**synopsis:** Famille : Groupe de personnes réunies par des liens de parenté et un fort sentiment de solidarité morale et matérielle. Quand Alain a épousé Nathalie, il ne savait pas qu'il épouserait aussi sa famille. Ce samedi, comme toutes les semaines, ils sont invités à dîner chez son beau-frère, Jean - Pierre à Créteil. Mais ce soir, plus que d'habitude, Alain est à bloc, il bout comme une cocotte prête à exploser. Il en a marre, marre de se planter à chaque fois sur le chemin pour aller à Créteil, marre de se taper les petits conseils de vie de Jean - Pierre et de sa femme Catherine qui élève ses enfants comme des **chevaux**, marre d'attendre de dîner l'estomac vide en regardant les spectacles soporifiques de leur fille Gaëlle, marre de regarder pour la énième fois la vidéo  
**titre:** Tellement proches  
**genre:** comédie

Par ailleurs, seulement les mots présent dans synopsis et le titre sont mis en surbrillance, (choix étant arbitraire). D'autres facettes que celles affichées sont calculées. Elles peuvent donc être affichées si on le souhaite en modifiant les facettes.

## III) Bibliographie :

[1]<http://vectors.nlpl.eu/repository/#> : fichier ayant l'id 51 pour le modèle de plongement de mots en français.

[2]<https://fasttext.cc/docs/en/crawl-vectors.html#models> : pour le modèle de plongement de mots en français.