

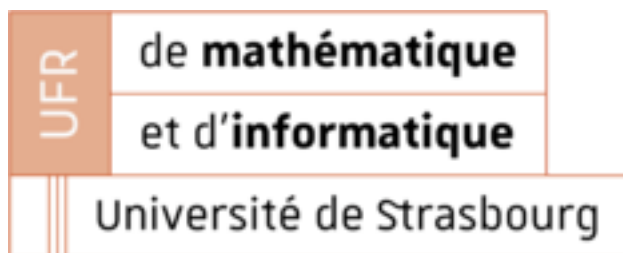
# Traitement Automatique de Langage

## Rapport de Projet

**Equipe:**

ALLEMAND Fabien

LEBOT Samuel



# Table des matières

Liste des figures	iii
<b>1 Analyse et Pré-traitement des Données</b>	<b>1</b>
<b>2 Méthodes Basiques</b>	<b>2</b>
2.1 Pré-traitement des Données . . . . .	2
2.2 Apprentissage Simple . . . . .	2
2.3 Apprentissage par Plis et Comparaison des Algorithmes . . . . .	3
<b>Bibliographie</b>	<b>4</b>

## Liste des figures

1	Effectifs des classes dans les données d'entraînement . . . . .	1
2	Comparaison des matrices de confusion pour des données déséquilibrées et équilibrées: les classes les moins représentées dans le jeu de données d'entraînement sont moins bien apprises lorsque les classes sont déséquilibrées. . . . .	2
3	Pipeline de pré-traitement des données . . . . .	2
4	Pipeline d'entraînement d'un modèle de régression logistique . . . . .	3
5	Précision de différents algorithmes lors d'un apprentissage par validation croisée avec 5 plis . . . . .	3

```
train_data["genre"].value_counts()

drame      501
comédie    483
romance     443
policier    331
horreur     299
science fiction 298
biopic      191
documentaire 167
historique  162
Name: genre, dtype: int64
```

Figure 1: Effectifs des classes dans les données d'entraînement

## Introduction

L'objectif du projet est de réaliser un système de recherche d'information dans une collection de descriptions de films publiées sur Allociné.

Le projet se décompose en deux parties:

- Prédiction du genre des films par TAL
- Visualisation des résultats

## 1 Analyse et Pré-traitement des Données

Dans un premier temps, les données d'entraînement (*allocine\_genres\_train.csv*) peuvent être chargées grâce à la fonction `read_csv` de la bibliothèque Pandas [2] en précisant le séparateur (`sep=","`). L'utilisation des méthodes `head`, `tail`, `describe`, `info` et `hist` permettent de visualiser et comprendre les données contenues dans le jeu de données complet.

Dans le cadre de ce projet seules les données contenues dans les colonnes *titre* et *synopsis* seront utilisées pour déduire la valeur contenue dans *genre*. Le jeu de données peut être réduit à ces trois features.

Etant donné que les données vont être utilisées pour de l'apprentissage automatique, il faut vérifier s'il y a des valeurs manquantes. Les méthodes `isna` et `sum` ne révèlent aucune valeur manquante dans les données d'entraînement.

La proportion des classes dans le jeu de données peut avoir un impact sur l'apprentissage. Les classes ayant un effectif plus faible seront généralement moins bien "appries". La figure 1 montre que les classes n'ont pas toutes la même proportion dans le jeu de données: il y a beaucoup d'individus de la catégorie *drame* alors que les classes *biopic*, *documentaire* et *historique* sont très peu représentées.

Après de nombreuses expériences, il s'avère que rééquilibrer les classes par *oversampling*, c'est à dire: dupliquer des individus des classes les moins représentées, donne de meilleurs résultats quelque soit la méthode utilisée. 2

**Dans tout la suite**, les résultats présentés correspondront aux résultats obtenus avec le jeu de données d'entraînement rééquilibré par *oversampling* à l'aide de l'objet `RandomOverSampler` de la bibliothèque Imbalanced-learn [1].

### Remarque

Le jeu de données d'entraînement contient trop peu de données pour effectuer un équilibrage des classes par *undersampling*, c'est à dire: supprimer des individus des classes les plus représentées. Les résultats obtenus avec cette méthode étaient généralement moins bons que sans équilibrage.

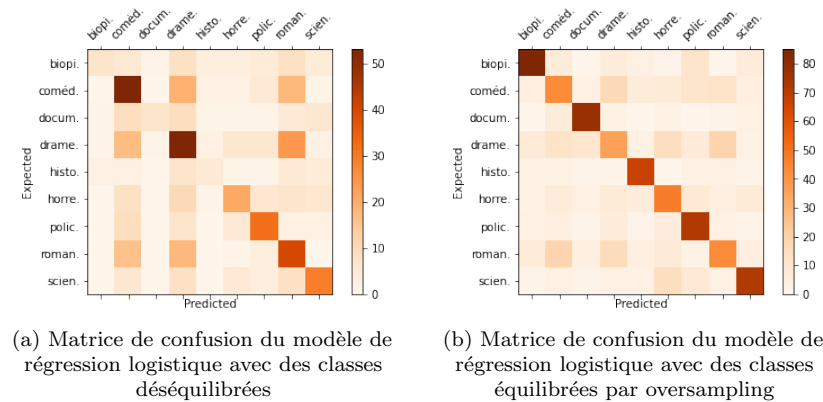


Figure 2: Comparaison des matrices de confusion pour des données déséquilibrées et équilibrées: les classes les moins représentées dans le jeu de données d'entraînement sont moins bien apprises lorsque les classes sont déséquilibrées.

```
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline

column_trans = ColumnTransformer(
    [
        # Titre: tf-idf
        ("titre_tfidf", titre_vectorizer, "titre"),
        # Synopsis: tf-idf
        ("synopsis_tfidf", synopsis_vectorizer, "synopsis"),
        # Synopsis: statistiques
        (
            "synopsis_stats",
            Pipeline(
                [
                    ("text_stats", text_stats_transformer),
                    ("vect", text_stats_vectorizer),
                    ("scaling", min_max_scaler)
                ]
            ),
            "synopsis"
        )
    ],
    # Others
    remainder="passthrough"
)
```

✓ 0.1s Python

Figure 3: Pipeline de pré-traitement des données

## 2 Méthodes Basiques

### 2.1 Pré-traitement des Données

Il est possible de mettre en place une pipeline 3 afin de vectoriser les données d'apprentissage et d'en extraire des informations statistiques.

Dans cette pipeline, les données déjà tokenisées vont être vectorisées selon la méthode TF-IDF en supprimant les stop-words. Puis des objets `FunctionTransformer` et `DictVectorizer` de la bibliothèque Scikit-learn [3] vont être utilisés pour obtenir les valeurs statistiques telles que la longueur du synopsis en nombre de mot et en nombre de phrases. Un `MinMaxScaler` est utilisé pour normaliser ces données afin qu'elles aient le même poids lors de l'apprentissage.

### 2.2 Apprentissage Simple

Avant tout apprentissage supervisé, il faut définir un jeu d'apprentissage et un jeu de test avec par exemple la méthode `train_test_split` de Scikit-Learn en spécifiant la proportion du jeu de données sélectionnée pour les données de test (typiquement: `test.size=0.2`) et précisant `shuffle=True` pour que les données ne soient pas

```
from sklearn.pipeline import make_pipeline
from sklearn.linear_model import LogisticRegression

# Pipeline pré-traitement et apprentissage
classifier_pipeline = make_pipeline(
    # Préparation des données pour l'apprentissage
    column_trans,
    # Algorithme d'apprentissage
    LogisticRegression()
)

# Apprentissage avec les données d'entraînement
classifier_pipeline.fit(X_train, y_train.to_numpy())
```

Figure 4: Pipeline d'entraînement d'un modèle de régression logistique

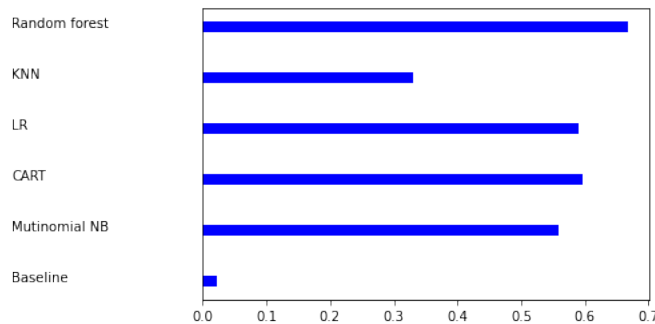


Figure 5: Précision de différents algorithmes lors d'un apprentissage par validation croisée avec 5 plis

sélectionnées séquentiellement (ce qui pourrait impacter le résultat si le jeu de données est trié par classes).  
**Remarque** || Le jeu de données *allocine\_genres\_test.csv* correspond au jeu de validation qui ne sera utilisé qu'après avoir définitivement choisi la méthode de prédiction.

Il est alors possible de créer une nouvelle pipeline pour automatiser le pré-traitement et l'apprentissage 4.

On peut ensuite faire des prédictions sur le jeu de test avec la méthode `predict` pour évaluer le modèle. Ce modèle de régression logistique donne une précision de 0.6 et un rappel de 0.61 (donc un score f1 de 0.6).

## 2.3 Apprentissage par Plis et Comparaison des Algorithmes

Comme vu précédemment, la taille et les individus sélectionnés pour entraîner le modèle peuvent avoir une influence sur la qualité de l'apprentissage. Pour contrer ce biais on peut utiliser des plis pour l'apprentissage: au lieu de diviser les données d'entraînement en un jeu d'entraînement et un jeu de test, on peut diviser les données d'entraînements en plusieurs plis qui seront tour à tour des données d'entraînement ou de test. Cela permet d'obtenir une idée de la performance moyenne du classifieur. On utilise la méthode *cross\_val\_score* de Scikit-Learn qui permet de conserver la proportion d'individus de chaque classe dans les plis.

En utilisant cette méthode, on peut comparer le modèle de régression logistique avec d'autres algorithmes. Les meilleurs résultats proviennent d'un modèle de forêt aléatoire (précision proche de 0.7). ??

## Conclusion

## Bibliographie

- [1] Imbalanced-learn. <https://imbalanced-learn.org/stable/>. Accessed: 2023-04-30.
- [2] Pandas. <https://pandas.pydata.org/>. Accessed: 2023-04-30.
- [3] Scikit-learn. <https://scikit-learn.org/stable/>. Accessed: 2023-04-30.