



École Supérieure de Génie Informatique

TD 1

Auteur:
AL ASMI Alhusain.

Cours:
Traitement Automatique de langue
Naturelle
Référent:
FOUCHET Arnaud

Contents

| | | |
|----------|---|----------|
| 1 | Introduction | 2 |
| 2 | Imports | 3 |
| 3 | Interface en Ligne de Commande (CLI) | 3 |
| 3.0.1 | Train | 3 |
| 3.1 | predict | 4 |
| 3.2 | evaluate | 4 |
| 4 | Résultats | 5 |
| 4.1 | Random Forest | 5 |
| 4.2 | Gradient Boosting Classifier | 5 |
| 4.3 | Ada Boost Classifier | 5 |

1 Introduction

Ce TD de Traitement Automatique du Langage Naturel (NLP), se penche sur un défi : la prédiction du caractère comique d'une vidéo en utilisant uniquement son titre.

Dans ce TD, nous explorons le potentiel des techniques NLP pour résoudre ce problème en utilisant Python et plusieurs bibliothèques spécialisées.

Le code que je présenterai dans ce rapport est un code modulaire et flexible permettra d'entraîner un modèle, de le sauvegarder, de réaliser des prédictions, et d'évaluer ses performances, le tout grâce à des commandes simples en ligne de commande.

2 Imports

Le code commence par plusieurs déclarations d'importation pour charger les bibliothèques et les fonctions nécessaires. Certaines de ces bibliothèques comprennent :

- **click** : Utilisé pour créer des interfaces en ligne de commande.
- **sklearn.model_selection** : Fournit des outils pour la sélection et l'évaluation de modèles, tels que la cross-validation.
- **pickle** : Utilisé pour sérialiser et désérialiser des objets Python, en particulier pour sauvegarder et charger des modèles entraînés
- **pandas** : Une bibliothèque populaire pour la manipulation des données et création des dataframes.
- Des modules personnalisés du package src, tels que `make_dataset`, `make_features`, `make_model` et `dump`.

3 Interface en Ligne de Commande (CLI)

La bibliothèque click est utilisée pour créer une interface en ligne de commande. Elle définit un groupe CLI nommé cli. Nous avons 3 commandes disponibles afin d'utiliser les 3 différentes fonctionnalités du code: **train**, **predict** et **evaluate**.

3.0.1 Train

Dans la fonction train, nous commençons par appeler la fonction `make_model()` pour créer un modèle et obtenir une configuration. Cette fonction lance un log dans la terminale permettant l'utilisateur de choisir les types de pré-traitement à appliquer sur les données.

```
Choose the configs you want to use:
1. Stemming
2. Stop words
3. TF-IDF
4. Default
5. Exit

Enter your choice: 1
Enter your choice: 2
Enter your choice: 5
(SignAIify) PS C:\Users\huso\PycharmProjects\NLP> 
```

Ces choix sont sauvegardés dans la liste de configurations et ensuite sauvegardée dans le même fichier Json que celui du modèle. Nous utilisons cette commande pour faire appel à la fonction train

```
python src/main.py train --task=is_comic_video --input_filename=src/data/raw/train.csv
--model_dump_filename=src/models/model.json
```

3.1 predict

La commande predict est définie avec des options pour la tâche, le nom du fichier d'entrée, le nom du fichier de sauvegarde du modèle et le nom du fichier de sortie. Cette commande est utilisée pour effectuer des prédictions avec le modèle entraîné.

Elle lit les données d'entrée à partir du fichier spécifié à l'aide de la fonction **make_dataset**. Ensuite, Elle charge un modèle pré-entraîné et une configuration à partir d'un fichier pickle sérialisé (model_dump_filename).

Similaire à la fonction train, elle applique les types de pré-traitement choisis et en suite les prédictions sont faites sur les titres. Une fois ceci est fait, les prédictions sont sauvegardées en tant que fichier CSV dans le répertoire "data/processed". Nous utilisons la commande suivante pour faire appel à la fonction predict.

```
python src/main.py predict --task=is_comic_video --input_filename=src/data/raw/test.csv  
--model_dump_filename=src/models/model.json --output_filename=src/data/processed/prediction.csv
```

3.2 evaluate

La fonction evaluate réutilise les mêmes fonctions déjà mentionnées, la différence dans cette fonction c'est qu'elle effectue la validation croisée sur le modèle à l'aide de cross_val_score et affiche les scores et la moyenne des scores. La commande permettant d'appeler la fonction evaluate est la suivante:

```
python src/main.py evaluate --task=is_comic_video --input_filename=src/data/raw/train.csv
```

4 Résultats

Après avoir testé 3 modèles en particuliers (Random Forest, Gradient Boosting Classifier, Ada Boost Classifier), et en utilisant différentes configurations de pré(traitement des données, chaque modèle a été soumis à une série de 8 tests distincts. Ces configurations comprenaient l'utilisation de `count_vectorizer`, `stemming`, `stop words`, et `TF-IDF`, ainsi que diverses combinaisons de ces techniques

Pour déterminer quels modèles sont meilleurs, nous pouvons examiner les moyennes des scores (Moyenne des scores) pour chaque type de modèle et configuration de pré-traitement des données. Plus la moyenne est élevée, meilleure est la performance du modèle. Voici une évaluation des meilleurs modèles en fonction des moyennes des scores :

4.1 Random Forest

- Configuration : `TF-IDF`
- Moyenne des scores : 0.9124

4.2 Gradient Boosting Classifier

- Configuration : `count_vectorizer`, `stemming`
- Moyenne des scores : 0.9249

4.3 Ada Boost Classifier

- Configuration : `count_vectorizer`, `stemming`
- Moyenne des scores : 0.9174

En fonction des moyennes des scores, le modèle Gradient Boosting Classifier avec la configuration `count_vectorizer` et `stemming` (Modèle 2) semble être le meilleur parmi les trois types de modèles testés. Il est suivi de près par le modèle Ada Boost Classifier avec la même configuration (Modèle 2). Le modèle Random Forest avec la configuration `TF-IDF` (Modèle 4) obtient également de bonnes performances.

Pour des informations plus détaillées sur les performances de chaque modèle, y compris les résultats de chacun des 8 tests, il suffit de consulter le répertoire "Resultats". Vous pouvez retrouver un fichier dédié à chaque modèle.