

# Détection automatique de la subjectivité dans les textes

BIENASSIS Jules, DELMOTE Adrien, MEUNIER Mortimer, MOUYSET Martin, OLEIWAN Jo  
Encadrant : MOTHE Josiane



## Introduction

### Subjectivité

- La subjectivité est vue comme une empreinte personnelle que le locuteur laisse dans son énoncé, reflétant ses positions, attitude et sentiments [1].
- La capacité à discerner si un texte est influencé par les préjugés de son auteur, soulignant la complexité de la tâche en raison de la variabilité des niveaux d’expertise et des biais personnels [2].
- L’interprétation individuelle du langage et des expériences personnelles [3].

### Utilisation

- La subjectivité permet l’analyse de sentiment [4].
- Les systèmes de recommandation [5].
- La désinformation [6].

## Données et prétraitement

### Collection de test

Les modèles d’apprentissage automatique ont besoin d’être utilisés dans des collections de tests en lien avec la subjectivité pour se situer en termes de performance [6].

- La réponse aux questions multi perspective MPQA et MPQA gold.
- Movie review. CheckThat! [1].

### Augmentation de données

Les modèles d’apprentissage automatique demandent une quantité de données importante. Il faut donc augmenter les données. Cela permet également de pallier des problèmes de déséquilibres de classes dans certains jeux de données.

- Traduction.
- Reformulation + Génération.

## References

- [1] Ruggeri et al: On the Definition of Prescriptive Annotation Guidelines for Language-Agnostic Subjectivity Detection (2023)
- [2] Antici et al : SubjectivITA: An Italian Corpus for Subjectivity Detection in Newspapers. Springer International Publishing, International Conference of the Cross-Language Evaluation Forum for European Languages, (2021)
- [3] Aker et al : A Corpus for Sentence-level Subjectivity Detection on English News Articles, 1–6 (2019)
- [4] Pang, B. et Lee, L.: A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts.(2004)
- [5] Chaturvedi et al: Distinguishing between facts and opinions for sentiment analysis: Survey and challenges (2018)
- [6] Galassi et al : Overview of the CLEF-2023 CheckThat! Lab, (2023)
- [7] Kaliyar et al : FakeBERT: Fake news detection in social media with a BERT-based deep learning approach, 11765–11788 (2021)

## Modèle de représentation de texte et langues

### Plongement de mots

Les plongements de mots capture la sémantique et la relations contextuelles des mots afin d’améliorer l’apprentissage des modèles de traitement du langage naturel [7].

- GloVe.
- Word2Vec.
- Vecteurs issus de modèles pré-entraînés

### BERT, LLM\* et variantes multilingues

- BERT Base.
- BERT Large.
- RoBERTa.
- XML-RoBERTa.

3 types de grands modèles de langue :

- Autoregressifs.
- Langage masqué.
- Langue encodeur-décodeur.

\* LLM = Large Language Model qui se traduit parGrand Modèle de Langue

## Approches non neuronales pour la détection de subjectivité

### Apprentissage automatique

- Bayésien Naïf.
- K-plus proche voisins.
- Machine à vecteurs de support (SVM).
- Arbre de décision.
- Forêts aléatoire.

L’extraction de caractéristiques linguistiques est essentielle pour l’utilisation de ces modèles.

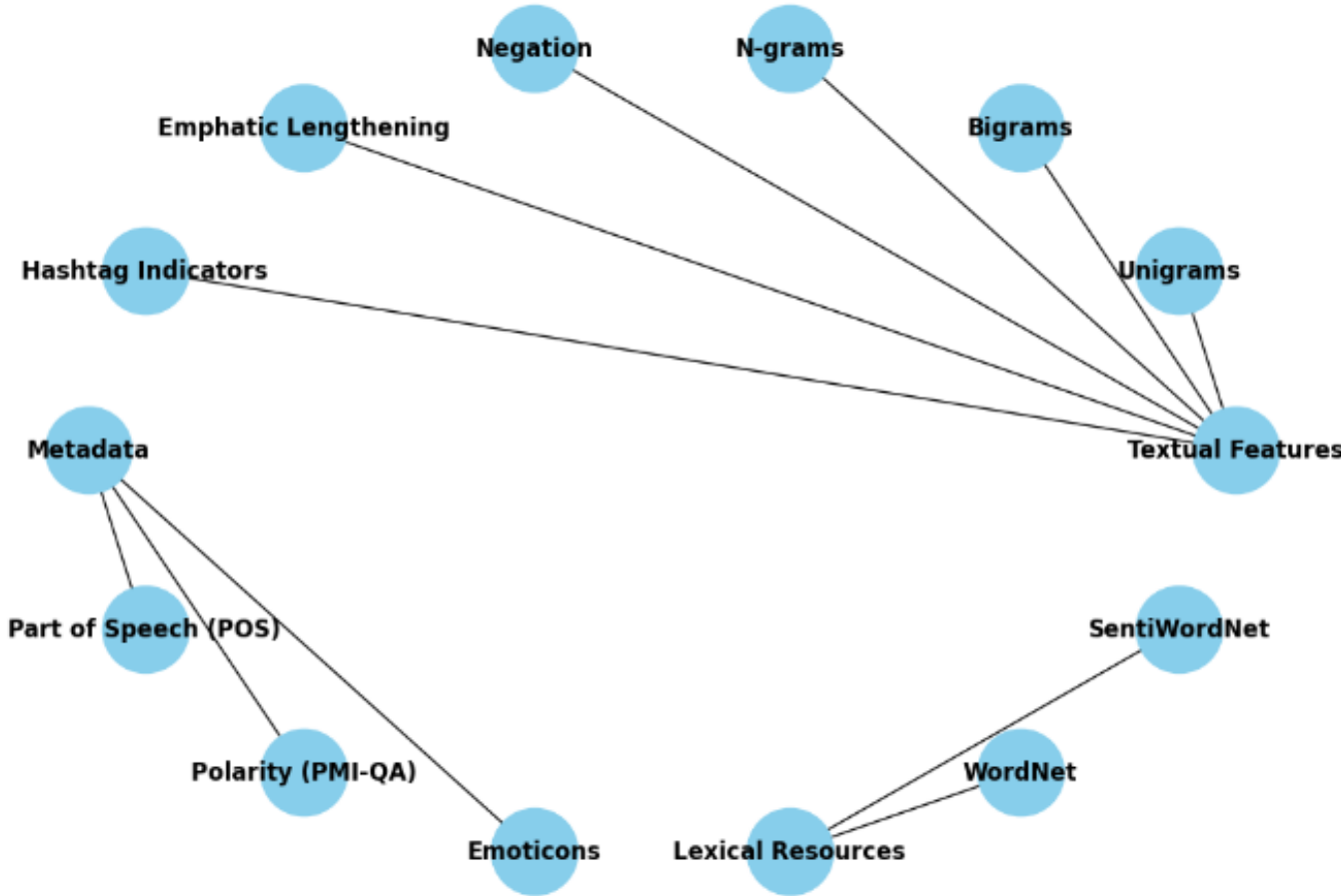


Figure 1. Familles de caractéristiques exploitables en NLP

### Graphe

Bo Pang et Lillian Lee [4] ont développées une technique qui permet de classer les phrases dont l’appartenance à une classe est trop floue. Ce modèle permet d’appliquer les algorithmes de gestion de flot maximum.

## Approches neuronales pour la détection de subjectivité

### CNN

Les réseaux de neurones convolutionnels (CNN) étaient considérés comme plus efficaces que les anciennes méthodes. Les filtres dans les couches supérieures du CNN peuvent capturer les relations syntaxiques entre les phrases éloignées dans la phrase en paramètre d’entrée. Certains modèles permettent également de prendre en compte le sentiment du texte, souvent ignoré dans les techniques précédentes.

### GRU

Gated Recurrent Unit (GRU) est une architecture basée sur les réseaux de neurones récurrents (RNN). Ces modèles utilisent une couche supplémentaire en sortie pour effectuer une classification de la subjectivité.

### FakeBERT

Avec une précision globale de 0.98, Kaliyar et al [7] ont montré qu’il est possible d’utiliser différentes couches de convolution entre un modèle de représentation d’une phrase, de type BERT ou GloVe appelé “Embedding”, et une couche dense devant effectuer la classification, comme visible sur la figure ci-contre.

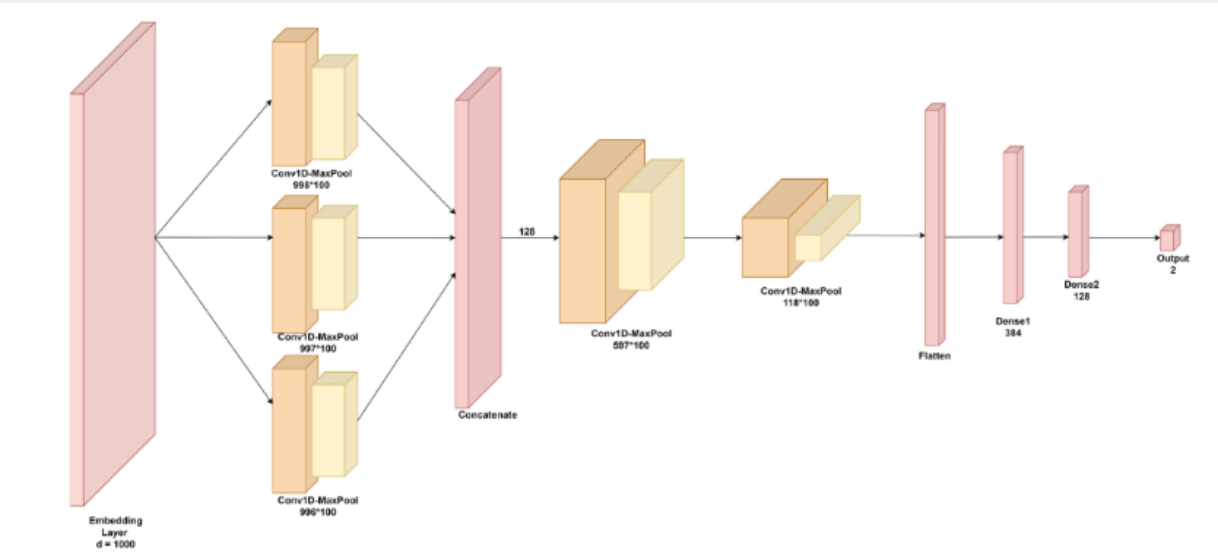


Figure 2. Architecture du modèle FakeBERT cf.[7]

### Ajustement fin de LLM

L’ajustement fin implique le réentraînement du modèle sur un ensemble de données plus petit et spécifique à une tâche auxiliaire pour avoir de meilleures performances sur celle-ci. Cet ajustement est établi à travers diverses méthodes.

- L’apprentissage basé sur les invites.
- L’apprentissage basé sur les patrons.
- L’apprentissage basé sur les tâches proxy.

## Limites et perspectives

### Limites

- Traduction: introduit biais et inexactitudes.
- GPT: Crée divergences stylistiques, affecte l’entraînement.
- Modèles ML: Hypothèses manuelles limitent la généralisation du modèle.
- CNN: Inefficaces pour dépendances à longue distance.
- GRU et LSTM: Plus de paramètres, apprentissage plus lourd, surpassés par les transformers.

### Perspectives

Utiliser BERT avec des couches comme CNN/LSTM peut améliorer la détection de subjectivité. La réutilisation et l’ajustement de modèles pré-entraînés favorisent les avancées en NLP, rendant la détection de subjectivité plus accessible et précise.