

Rapport Synthétique : Détection de Fausses Nouvelles par Apprentissage Profond

Benatmane Hamza

1 Données et Problématique

1.1 Jeu de Données

- **Source** : Dataset de détection de fausses nouvelles
- **Volume** : 44,898 articles (23,481 faux, 21,417 vrais)
- **Distribution** : Relativement équilibrée (52.3% faux, 47.7% vrais)
- **Caractéristiques** : Textes d'actualités de longueurs variables (moyenne $\sim 2,500$ caractères)

1.2 Problématique

Développement d'un système de classification binaire capable de distinguer les vraies des fausses nouvelles en utilisant un réseau de neurones multicouche (MLP).

2 Méthodologie

2.1 Prétraitement des Données

- Nettoyage du texte : suppression des URLs, caractères spéciaux, chiffres
- Normalisation : mise en minuscules, lemmatisation
- Vectorisation : TF-IDF avec 5,000 caractéristiques
- Division des données : 60% entraînement, 20% validation, 20% test

2.2 Architecture du Modèle

- **Réseau multicouche profond** :
 - Couche d'entrée : 5,000 neurones (dimensions TF-IDF)
 - Couches cachées : $512 \rightarrow 256 \rightarrow 64$ neurones

- Couche de sortie : 1 neurone (sigmoid)
- **Techniques de régularisation :**
 - Dropout (0.3, 0.2, 0.1)
 - Batch Normalization
 - Early Stopping

3 Résultats et Analyse

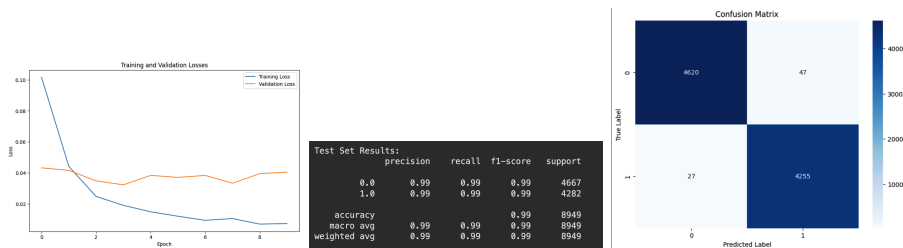


Figure 1: Performances du modèle

4 Pistes d'Amélioration

4.1 Améliorations Potentielles

1. Prétraitement

- Utilisation de techniques de data augmentation
- Exploration d'autres méthodes de vectorisation (Word2Vec, BERT)

2. Entraînement

- Implémentation de techniques d'apprentissage par transfert
- Utilisation de techniques d'ensemble

5 Conclusion

Le modèle développé montre des performances prometteuses pour la détection de fausses nouvelles. Les techniques de régularisation et l'architecture choisie permettent une bonne généralisation, mais il existe encore des pistes d'amélioration significatives, notamment dans l'utilisation de modèles plus avancés et de techniques de traitement du langage naturel plus sophistiquées.