

# Stylo@DEFT2022

Notation automatique de copies d'étudiant.e-s par  
combinaisons de méthodes de similarité

---

Ibtihel Ben Ltaifa<sup>1</sup>, Toufik Boubehziz<sup>1</sup>, Andrea Briglia<sup>1</sup>,  
Corina Chutaux<sup>1</sup>, Yoann Dupont<sup>2</sup>, Carlos-Emiliano González-Gallardo<sup>3</sup>,  
Caroline Koudoro-Parfait<sup>1,2,4</sup> Gaël Lejeune<sup>1,3</sup>

July 5, 2022

<sup>1</sup> Sens Texte Informatique Histoire (STIH), Sorbonne Université

<sup>2</sup> Observatoire des Textes et des Connaissances (OBTIC), Sorbonne Université

<sup>3</sup> Laboratoire Informatique, Image et Interaction (L3i), La Rochelle Université, France

<sup>4</sup> Sorbonne Center for Artificial Intelligence (SCAI), Sorbonne Université, France

1. Méthodes de similarité pour la correction automatique
2. Protocole expérimental
3. Résultats et discussion
4. Conclusion et perspectives

Évaluation automatique de devoirs d'étudiant-e-s avec le corrigé de l'enseignant



**Figure 1:** Comment faire pour ne pas mettre que des 1 et des 0 ? Source : Français Langue Seconde <sup>2</sup>

<sup>1</sup><http://www.francaislangueseconde.fr>

<sup>2</sup><http://www.francaislangueseconde.fr>

Deux ressources :

- questions (50 train et 21 test)+ éléments de réponse (et indication de notation)
- réponses (3.820 train et 1.644 test) produites par des étudiants (avec la note)

	Nb. Questions	Nb. Réponses	Nb. Etudiants
Données d'entraînement	50	3 820	118
Données d'évaluation	21	1 644	118

**Table 1:** Données fournies pour la tâche de notation automatique

# Méthodes de similarité pour la correction automatique

---

→ Approche basée sur l'usage de plusieurs méthodes de similarité pour la correction automatique

1. Extraction de caractéristiques ;
2. Mesures de similarité

→ Quatre types de caractéristiques pour la vectorisation des réponses d'étudiant.e.s

1. N-grammes sans frontières de mots (char) ;
2. N-grammes avec frontières de mots (char\_wb);
3. Sous-chaînes de mots (WordPiece) ;
4. Vecteurs contextualisés SentenceBERT (sbert)

→ Deux mesures de similarité sont appliquées sur les différentes représentations vectorielles des textes générées à partir de quatre différents types de caractéristiques:

1. La similarité Cosinus ;
2. L'indice de dissimilarité de Bray-Curtis



# Protocole expérimental

---

Nous utilisons trois méthodes de notation différentes :

1. notation par régression;
2. notation par classification;
3. notation par réseau de neurones

# Notation par régression (M1)

→ tâche de notation définie comme une tâche de régression linéaire.

Nous avons utilisé les traits suivants :

1. char ;
2. char\_wb ;
3. WordPiece ;
4. sbert

## Notation par classification (M2)

- tâche de notation définie comme une tâche de classification avec 3 classes (0 , 0,5 , 1) en utilisant:
- un algorithme de régression logistique avec le solveur Saga (meilleurs résultats).
- les mêmes traits que M1

## Notation par réseau de neurones (M3)

- Appel à un modèle de réseau de neurones pour effectuer la notation.
- Le modèle est constitué de deux couches de neurones associées à une fonction d'activation de type sigmoïde.
- L'algorithme prend en entrée les mêmes traits de M1 et la note prédite en sortie.

## Résultats et discussion

---

	M1	M2	M3
Précision (P)	0,642	0,678	0,653
Corrélation de Spearman ( $r_s$ )	0,633	0,543	0,644

**Table 2:** : Résultats de nos méthodes entraînées sur DEFT-TRAIN-2021 et évaluées DEFT-TEST2021

Représentation	Longueurs	Distances	# Dimensions
char	[3 : 3], [3 : 4],[3 : 5] . . . [4 : 6]	Cosinus et Bray-Curtis	16
char_wb	[3 : 3], [3 : 4],[3 : 5] . . . [4 : 6]	Cosinus et Bray-Curtis	16
WordPiece			16

**Table 3:** : Les 48 dimensions utilisés pour la représentation de la similarité (nous y ajoutons deux représentations issues de Sentencebert pour aboutir à 50 dimensions)



# Analyse de nos résultats

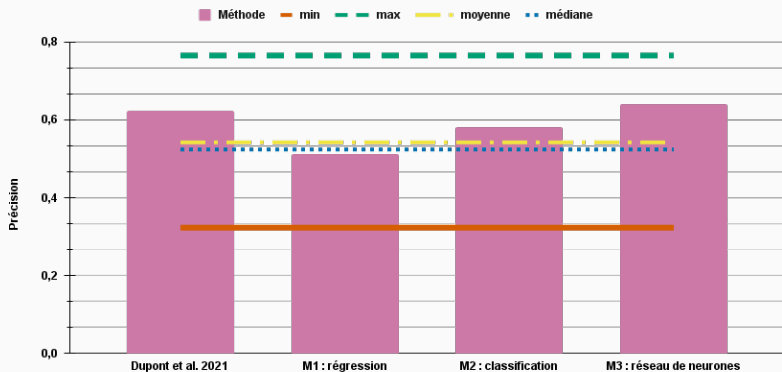


Figure 2: Résultats des nos méthodes et statistiques générales

- Les résultats obtenus sont supérieurs à la moyenne et à la médiane des résultats des autres équipes, méthode avec un certain potentiel
- améliorer la précision de **0,630** à une précision de **0,678** (M3), soit un gain de **4,8** points de pourcentage sur DEFT 2021.
- score de **0,654** pour DEFT 2022, soit un gain de **2,4** points de pourcentage par rapport à DEFT 2021.

## **Conclusion et perspectives**

---

## Proposition d'une approche simple pour l'évaluation de réponses d'étudiant-e-s étant donné une correction préalable :

- Nous avons pu améliorer les résultats depuis notre précédente contribution (DEFT2021)
- Notre score a augmenté de **4,8** points de pourcentage sur les données Deft 2021, **2,4** point de pourcentage sur les données Deft 2022.

- Tester d'autres similarités comme la distance de Jaccard ou l'indice de Dice pour estimer la proximité.
- Utiliser d'autres modèles de sentence embedding comme le modèle : Universal Sentence Encoder (Cer et al., 2018 ; Yang et al., 2019)