

# Classification multi-classes de descriptions d'emploi

Pol Labarbarie et Mehdy Hounkonnou

Université de Bordeaux

Vendredi 15 Janvier 2021

## Qu'est-ce que le NLP ?

- 1 NLP (Natural language processing) est un sous-domaine de l'informatique, de la linguistique et de l'IA.
- 2 Compréhension du contenu des documents pour diverses applications dont la classification.
- 3 Exemples : Siri (Apple), Google Translate, Alexa (Amazon).

## Le Défi IA

- 1 Classification multi-classes avec 28 classes au total.
- 2 Données extraites de CommonCrawl, formant le modèle GPT-3 d'OpenAI à priori biaisé.
- 3 Objectif : Concevoir une solution qui est à la fois précise et équitable.

## Conditions

Tous les temps de calcul sont obtenus avec un processeur Intel-Core i7-7820HQ (8 coeurs 2.9 GHz, 3.9GHz Turbo) aka "The quiet BEAST".

- 1 Statistiques descriptives et pré-traitement
- 2 Amélioration de la Baseline
- 3 Transformers et réseaux de neurones profonds
- 4 Conclusion

## Descriptif du jeu de données

- 3 jeux de données : apprentissage, test et label.
- Nous disposons de descriptions d'emploi et l'objectif est de prédire un métier.
- Le jeu de données apprentissage :  $n = 217197$  lignes (descriptions).
- Le jeu de données test :  $n = 54300$  descriptions.
- Cadre d'apprentissage supervisé.

## Caractéristiques du jeu de données

- Aucune donnée manquante.
- Présence de langues autres que l'anglais :

03 - "Dr. Maria Ignasia Tjahjadi practices at Rumah Sakit Ibu dan Anak Hermina in Tanjung Priok, Jakarta. She completed S.Ked. from Atma Jaya University in 1996."

- "Soy periodista con 10 años de experiencia trabajando en medios digitales. He sido redactora web y coordinadora de redacción. Tengo conocimientos en lenguaje..."

- Au total : 93 descriptions avec des mots non anglais.

## Vocabulaire du jeu de données

- Présence de chiffres.
- Présence de mots sans signification dans le corpus.
- Présence de combinaisons de mots et de chiffres dépourvu de sens lexical ou de symboles.  
Exemple : '2006music', '.\_.', 'aaaahc'

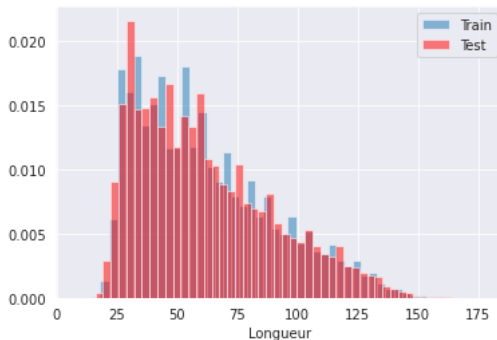


Figure: Histogrammes des données brutes

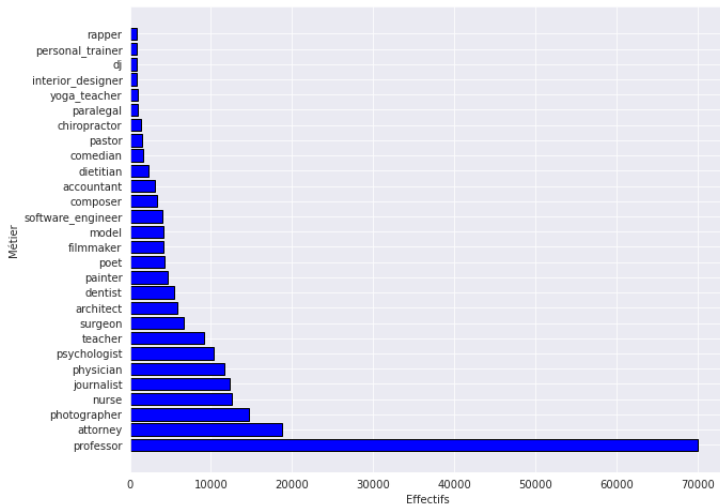


Figure: Histogrammes des répartitions des emplois

## Traitement

- ❶ Mise en minuscule des descriptions.
- ❷ Création d'un dictionnaire des abréviations pour les étendre dans les descriptions.  
Exemple: ("gf" = "girlfriend", "cuz" = "because", "coz" = "because")
- ❸ Suppression des caractères spéciaux.
- ❹ Lemmatisation.

	original_word	lemmatized_word
0	trouble	trouble
1	troubling	trouble
2	troubled	trouble
3	troubles	trouble
	original_word	lemmatized_word
0	goose	goose
1	geese	goose

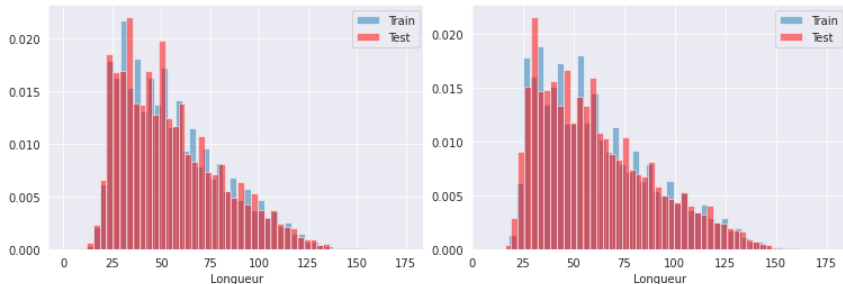
Figure: Exemples du processus de lemmatisation

- ❺ Suppression des mots non-anglais.

## Traitement

1640173 mots non anglais enlevés.

Nous remarquons clairement une réduction de la longueur des phrases.



(a) Histogramme des données lemmatisées

(b) Histogramme des données brutes

Figure: Comparaison des histogrammes des longueurs des descriptions



① Statistiques descriptives et pré-traitement

② Amélioration de la Baseline

③ Transformers et réseaux de neurones profonds

④ Conclusion

## La Baseline

- ① Pas de traitement des descriptions.
- ② Pondération TF-IDF appliquée au jeu de données.
- ③ Régression logistique avec pénalité Ridge.

Résultat : Taux de bonne classification 0.72.

## TF-IDF qu'est-ce que c'est ?

TF-IDF (term frequency-inverse document frequency) est une mesure statistique qui évalue la pertinence d'un mot pour un document appartenant à une collection de documents.

## TF-IDF Mathématiquement

Pour le mot  $i$  dans le document  $j$ , on pose  $tf_{i,j} = \frac{n_{i,j}}{\sum_n n_{i,j}}$

$$w_{i,j} = tf_{i,j} \log\left(\frac{N}{df_i}\right)$$

$n_{i,j}$  nombre d'apparitions du mot  $i$  dans le documents  $j$ ,

$N$  le nombre total de documents,

$df_i$  le nombre de documents contenant le mot.

## Résultats

La baseline donne un score sur Kaggle de 0.73244. Les résultats en local sont affichés à droite.

```
Accuracy: 0.78
Auc: 0.98
Detail:
```

	precision	recall	f1-score	support
0	0.66	0.42	0.52	293
1	0.79	0.74	0.76	818
2	0.75	0.57	0.65	182
3	0.53	0.53	0.53	1881
4	0.80	0.62	0.69	146
5	0.81	0.78	0.79	921
6	0.65	0.75	0.69	2516
7	0.83	0.49	0.62	188
8	0.82	0.68	0.74	1330
9	0.81	0.65	0.72	632
10	0.77	0.59	0.67	161
11	0.73	0.69	0.71	2339
12	0.84	0.66	0.74	342
13	0.69	0.61	0.65	797
14	0.87	0.80	0.83	2532
15	0.75	0.66	0.70	845
16	0.94	0.88	0.91	1096
17	0.88	0.55	0.67	284
18	0.79	0.74	0.77	775
19	0.81	0.89	0.85	13898
20	0.81	0.88	0.84	2878
21	0.84	0.60	0.70	152
22	0.76	0.67	0.71	2151
23	0.92	0.27	0.42	209
24	0.74	0.59	0.66	1162
25	0.80	0.80	0.80	670
26	0.83	0.87	0.85	3779
27	0.88	0.76	0.81	463
accuracy			0.78	43440
macro avg	0.79	0.67	0.71	43440
weighted avg	0.78	0.78	0.78	43440

Figure: Résultats de la baseline

## Résultats

La matrice de confusion de la baseline est la suivante. On note qu'il se trompe beaucoup entre les professeurs, les teachers (enseignants) et les psychologues (classe 19, 3 et 22).

		Confusion matrix																												
True	0	124	2	0	54	1	0	22	0	1	1	0	1	1	1	4	2	0	0	2	40	2	1	6	0	3	3	22	0	
	1	0	608	0	33	1	2	55	1	0	4	0	2	9	2	9	1	1	0	4	14	46	0	3	0	4	1	16	2	
	2	0	4	0	104	39	2	1	6	0	0	1	1	1	0	0	3	0	0	1	2	3	0	9	0	0	1	1	3	
	3	19	27	21	1001	1	17	114	1	1	3	3	4	3	9	18	35	0	1	8	418	46	0	58	0	5	17	49	2	
	4	0	5	1	10	90	0	7	0	0	0	0	5	0	0	3	0	0	1	0	6	3	0	4	0	2	0	5	4	
	5	0	2	0	16	0	714	13	1	0	0	0	3	0	1	0	5	0	0	5	48	94	0	3	0	7	4	5	0	
	6	3	12	0	80	0	8	1878	1	2	4	1	6	6	3	4	25	0	0	26	281	84	3	26	0	10	2	50	1	
	7	0	5	0	9	0	2	13	93	0	0	0	0	0	0	1	0	0	0	14	18	0	3	0	23	0	7	0	0	
	8	3	7	0	13	1	0	14	0	903	1	0	142	0	1	8	2	34	1	2	171	8	0	2	1	1	1	13	1	
	9	3	5	0	26	0	1	36	0	0	411	2	1	0	2	3	3	1	0	1	60	8	0	4	0	11	1	51	2	
	10	0	4	0	6	0	0	18	0	0	0	95	0	3	3	1	0	0	0	2	3	6	7	0	0	1	10	2	0	
	11	0	1	2	22	3	1	43	0	77	3	0	1612	0	2	88	8	10	13	1	386	9	0	27	1	0	1	25	4	
	12	2	17	0	17	0	0	35	0	0	0	1	1	226	0	0	0	0	0	4	10	9	2	5	0	2	2	9	0	
	13	1	1	0	24	0	0	42	0	0	2	0	2	2	486	2	4	0	0	2	140	18	0	3	0	62	2	4	0	
	14	3	9	2	58	2	1	30	0	5	2	0	165	0	0	2024	3	0	0	2	134	17	0	45	1	1	1	20	7	
	15	1	2	0	41	0	9	48	0	0	0	2	0	4	1	1	560	0	0	9	118	26	1	3	0	2	12	5	0	
	16	1	1	0	13	0	0	9	0	13	1	0	27	0	0	8	0	961	2	0	35	8	0	7	0	0	1	9	0	
	17	0	0	1	2	2	0	4	0	5	1	0	70	0	0	6	0	4	155	0	15	2	0	6	0	0	0	7	4	
	18	1	5	0	15	0	9	49	0	0	0	1	0	6	1	0	5	0	0	575	54	43	0	2	0	0	5	4	0	
	19	16	4	1	190	4	48	179	3	89	34	1	137	0	46	80	62	9	1	43	12365	57	1	198	0	80	46	186	14	
	20	1	21	0	28	0	50	85	2	0	1	2	0	3	5	1	2	0	0	21	71	2541	0	8	0	12	4	18	2	
	21	0	9	0	3	0	1	11	0	0	0	9	1	4	0	0	5	0	0	0	3	3	91	0	0	2	9	1	0	
	22	2	4	7	80	2	4	56	0	9	3	0	23	0	3	18	17	3	1	1	418	24	0	1433	0	7	0	35	1	
	23	0	1	0	3	1	0	8	0	0	6	0	1	0	0	1	0	0	0	0	15	2	0	2	57	1	0	111	0	
	24	1	2	0	28	0	12	21	10	2	4	0	1	0	135	1	3	0	0	4	183	30	0	10	0	691	3	21	0	
	25	1	4	0	22	0	1	14	0	0	0	6	1	0	0	0	5	0	0	6	56	12	2	0	0	2	535	3	0	
	26	5	9	0	37	0	1	87	0	0	24	0	5	2	6	18	4	0	1	8	241	17	0	15	2	6	4	3285	2	
	27	0	1	0	12	3	0	12	0	0	1	0	9	0	1	16	0	0	0	0	37	5	0	12	0	1	0	2	351	
			0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27

Figure: Matrice de confusion de la baseline

# Comment améliorer la Baseline ?

## Réduction de dimensionalité

L'ACP ne conduisant pas à de bons résultats. → Sélection de variables par test du khi-deux.

→ Permet d'obtenir un modèle plus parcimonieux mais plus interprétable que pour une ACP.

## Agréger des prédicteurs réduit la variance

Si on note pour tout  $1 \leq l \leq q$   $Z_l$  un prédicteur, on peut montrer que

$$\text{Var} \left( \frac{1}{q} \sum_{l=1}^q Z_l \right) = \sigma^2 \frac{1-\rho}{q} + \sigma^2 \rho \leq \sigma^2$$

où  $\text{Var}(Z_l) = \sigma^2$  et  $\text{corr}(Z_l, Z_{l'}) = \rho$

## Choix du prédicteur

Machine à vecteurs de support (SVM) conduit à des meilleurs résultats (quelques %) que la Régression logistique.

→ Utilisation de Machine à vecteurs de support.

## Méthodologie

- 1 TF-IDF sans stop-words et avec notre corpus traité.
- 2 Choix du nombre de variables conservées du test du khi-deux obtenu par validation croisée 5-folds.
- 3 Agrégation (Bagging) de 100 Machines à vecteurs de support.

## Résultats de la baseline améliorée

Les résultats de la baseline améliorée en local sont les suivantes, donnant un score sur Kaggle de 0.74231.

	precision	recall	f1-score	support
0	0.64	0.49	0.56	293
1	0.75	0.74	0.74	818
2	0.67	0.65	0.66	182
3	0.55	0.49	0.52	1881
4	0.68	0.64	0.66	146
5	0.77	0.79	0.78	921
6	0.66	0.73	0.69	2516
7	0.77	0.54	0.64	188
8	0.78	0.70	0.73	1330
9	0.75	0.66	0.70	632
10	0.76	0.60	0.67	161
11	0.71	0.69	0.70	2339
12	0.82	0.68	0.74	342
13	0.65	0.58	0.61	797
14	0.86	0.79	0.83	2532
15	0.71	0.67	0.69	845
16	0.91	0.90	0.90	1096
17	0.85	0.61	0.71	284
18	0.74	0.76	0.75	775
19	0.81	0.88	0.85	13898
20	0.82	0.86	0.84	2878
21	0.76	0.64	0.69	152
22	0.76	0.67	0.71	2151
23	0.77	0.30	0.43	209
24	0.70	0.58	0.64	1162
25	0.80	0.78	0.79	670
26	0.84	0.86	0.85	3779
27	0.82	0.81	0.82	463
accuracy			0.78	43440
macro avg	0.75	0.68	0.71	43440
weighted avg	0.78	0.78	0.77	43440

Figure: Résultats de la baseline améliorée

## Résultats

La matrice de confusion de la baseline améliorée est la suivante. On note à une légère amélioration près les mêmes confusions de classes(19,3, 22..).

Matrice de confusion																												
0	144	6	0	30	1	0	20	0	1	1	0	2	2	3	4	4	0	0	2	45	1	1	5	1	3	4	13	0
1	1	606	1	30	4	2	43	0	3	5	1	3	6	5	9	1	1	0	8	22	39	2	2	0	6	3	12	3
2	0	1	118	22	4	0	6	0	1	2	0	1	0	0	5	0	0	2	2	5	1	6	0	0	1	1	1	2
3	19	31	36	920	3	22	117	1	4	7	2	10	6	11	16	43	6	1	12	445	40	0	50	1	8	16	49	5
4	0	3	2	7	93	0	3	0	1	1	0	6	0	0	3	1	0	2	0	8	1	0	2	0	1	0	3	9
5	3	6	0	16	0	725	10	2	2	1	0	2	0	2	0	3	0	0	4	50	73	0	4	0	8	6	4	0
6	5	17	1	73	1	9	1828	1	1	6	3	10	14	6	9	33	0	0	32	284	77	5	29	1	16	3	48	4
7	0	3	0	6	0	2	9	102	0	0	0	4	2	0	1	0	0	1	17	13	0	4	0	21	0	3	0	
8	3	7	0	10	2	1	16	0	929	2	0	126	1	0	22	2	40	1	4	132	8	1	5	0	3	1	12	2
9	4	7	0	23	0	2	29	0	3	418	0	4	1	5	2	2	2	0	0	54	8	0	5	3	8	2	48	2
10	0	5	0	5	0	1	14	0	0	1	96	0	2	3	0	1	0	0	4	4	4	7	1	0	1	10	2	0
11	3	6	1	20	3	3	39	0	92	3	0	1617	2	1	77	7	11	13	0	356	8	0	33	1	2	2	27	12
12	1	15	0	17	0	0	28	0	1	0	0	1	232	3	0	2	0	0	6	11	3	3	6	0	1	2	10	0
13	0	5	0	29	0	0	36	0	2	4	1	4	2	460	3	5	0	0	3	157	11	0	5	1	59	3	7	0
14	3	10	3	50	4	3	34	1	17	2	0	178	2	1	1999	5	3	1	3	127	18	0	39	0	3	1	17	8
15	0	2	0	48	0	11	39	1	0	1	3	3	2	3	1	563	0	0	11	102	22	0	8	0	7	11	7	0
16	2	1	0	9	1	0	7	1	16	2	0	29	0	0	5	0	981	2	0	24	7	0	4	0	1	1	3	0
17	0	1	1	0	2	0	2	0	2	1	0	63	0	0	7	0	3	174	1	11	1	0	4	0	1	0	6	4
18	0	8	1	8	0	10	31	0	0	2	2	3	4	2	3	4	0	0	590	43	44	0	4	0	4	4	8	0
19	19	2	1	166	7	59	185	7	106	42	1	154	3	47	94	70	13	4	53	1227	55	4	196	1	92	41	175	24
20	1	27	0	22	2	64	105	3	3	4	1	3	0	7	9	3	0	26	81	2470	2	9	0	19	5	12	0	
21	0	8	0	3	0	2	5	0	0	1	7	2	1	1	0	3	1	0	0	5	2	97	0	0	3	7	3	1
22	5	7	10	74	1	8	50	0	4	4	1	15	1	8	21	16	10	3	2	410	26	1	1437	0	4	2	29	2
23	0	1	0	5	1	0	5	0	0	3	0	2	0	0	0	1	0	0	0	13	1	0	3	62	0	0	112	0
24	1	6	0	22	0	12	26	14	5	8	0	0	0	136	0	4	2	0	6	167	45	1	8	0	677	5	17	0
25	2	7	0	22	0	1	11	0	0	1	9	5	0	1	0	14	1	0	11	44	9	3	1	0	3	523	2	0
26	8	12	0	48	1	4	83	0	4	32	0	11	1	7	9	8	2	1	13	232	21	0	10	10	10	1	3249	2
27	0	3	0	3	7	0	6	0	1	1	0	10	0	0	12	0	0	0	0	29	5	0	7	0	1	0	3	375
	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27

Figure: Matrice de confusion de la baseline améliorée

## Essai sans succès

Même méthodologie que précédemment mais en unifiant les données. → Pas de bons résultats.

## Améliorations

- Le test du khi-deux est effectué sur 40000 variables, il faudrait utiliser des méthodes de multiplicité des tests (ex : Bonferroni).
- L'agrégation de prédicteurs permet d'obtenir de meilleurs résultats mais n'améliore pas significativement les résultats, nous ne faisons que réduire la variance.
- D'après (Joulin et al. 2016), les méthodes standards comme les SVM et les RL sont limités en terme de performance pour la classification de document.

## Solution

→ Les réseaux de neurones profonds.



- 1 Statistiques descriptives et pré-traitement
- 2 Amélioration de la Baseline
- 3 Transformers et réseaux de neurones profonds**
- 4 Conclusion

## L'état de l'art

D'après (Bojanowski 2019; Neveu n.d.), les méthodes d'apprentissage profond sont l'état de l'art pour la classification de document.

## Plongement lexical (Word-embedding)

Associe à chaque mot d'un vocabulaire un vecteur de nombre réels généralement de taille 300.  
→ Cette représentation permet de capter la sémantique des mots.

- Word2Vec (Mikolov et al. 2013)
- GloVe (Pennington and Socher 2014)

- Modèle entraîné à : Selon un mot prédire  $c$  mots de contexte autour de lui.
- Fonction softmax en sortie.
- Par rétropropagation du gradient on obtient les formules de mise à jour des poids de la couche  $W$  et  $W'$ .

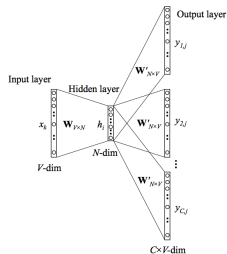


Figure: Architecture du modèle Skip-Gram

Utilisation d'une représentation vectorielle précise → Téléchargement des vecteurs de mot.

## Méthodologie

Utilisation du modèle de l'article vu en cours (Pietro n.d.).

- ① Conversion de notre jeu de données apprentissage et test en index se référant au vocabulaire de la représentation vectorielle.
- ② Entraînement d'un réseau avec deux couches LSTM sur 15 epochs.

## Résultats

- 1 epoch = 45min - 1h.
- Taux de bons classements sur les données Kaggle d'environ 0.71.
- Difficultés à la généralisation → Transformer.

## Les Transformers c'est quoi ?

C'est un empilement de couches aux architectures identiques. Chaque couche est composée d'une couche d'attention et d'une couche à propagation avant (Bojanowski 2019; Team n.d.).

DistilBERT est une version distillée de BERT. Plus rapide, plus économique.

## Méthodologie

- ① Utilisation des descriptions non pré-traitées.
- ② Tokenization avec la propre fonction du réseau.
- ③ Chaque séquence (description tokénisée) est de longueur maximale 80. (padding, truncation)
- ④ Réglage fin (fine-tuning) du modèle sur 3 epochs sinon sur-apprentissage.

## Résultats

- 1 epoch = 2h.
- Mieux que la Baseline mais pas mieux que notre Baseline améliorée  
→ Nécessité d'utiliser un Transformer plus gros → Serveur de calcul.

- 1 Statistiques descriptives et pré-traitement
- 2 Amélioration de la Baseline
- 3 Transformers et réseaux de neurones profonds
- 4 Conclusion

## C'est déjà la fin

- Les méthodes standards ne sont pas très performantes mais permettent d'obtenir une Baseline solide.
- Pour des résultats meilleurs, il faut utiliser un réseau Transformer.
- Plus le réseau est gros, meilleurs seront les résultats.
- Utilisation d'un autre réseau Transformer pour synthétiser des descriptions en plus pour rééquilibrer les classes.



Pitor Bojanowski. *Le langage naturel*. <https://www.college-de-france.fr/site/stephane-mallat/seminar-2019-02-20-11h15.htm>. Février 2019.



Armand Joulin et al. “Bag of Tricks for Efficient Text Classification”. In: *arXiv preprint arXiv:1607.01759* (2016).



Tomas Mikolov et al. “Efficient estimation of word representations in vector space”. In: *arXiv preprint arXiv:1301.3781* (2013).



Thibault Neveu. *L'état de l'art du NLP - Julien Chaumond, CTO Hugging Face - Podcast IA*. <https://www.youtube.com/watch?v=RL8QQk-LJp8t=314s>.



Jeffrey Pennington and Socher. “Glove: Global vectors for word representation”. In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014, pp. 1532–1543.



Mauro Di Pietro. *BERT for Text Classification with NO model training*. <https://towardsdatascience.com/text-classification-with-no-model-training-935fe0e42180>.



The Hugging Face Team. *Librairie Transformers Python*. <https://huggingface.co/transformers/quicktour.html>.