



HAL
open science

Défi TextMine 2024: "Reconnaissance d'entités géographiques dans un corpus des Instructions nautiques" - soumission équipe CRIT

Nicolas Gutehrlé

► To cite this version:

Nicolas Gutehrlé. Défi TextMine 2024: "Reconnaissance d'entités géographiques dans un corpus des Instructions nautiques" - soumission équipe CRIT. TextMine'24, Jan 2024, Dijon, France. hal-04442175

HAL Id: hal-04442175

<https://hal.science/hal-04442175v1>

Submitted on 6 Feb 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Défi TextMine 2024 : "Reconnaissance d'entités géographiques dans un corpus des Instructions nautiques" - soumission équipe CRIT

Nicolas Gutehrle*

*Université de Franche-Comté, CRIT
F-25000 Besançon, France
nicolas.gutehrle@univ-fcomte.fr,
<https://nicolasgutehrle.github.io/>

Résumé. Cet article présente notre participation au défi TextMine 2024. Nous proposons une méthode reposant sur un modèle CRF linéaire, ainsi qu'une méthode hybride combinant un modèle CRF linéaire et un modèle XGBoost. Nous entraînons les modèles sur des caractéristiques morphologiques et syntaxiques extraites de chaque mot. Notre modèle CRF linéaire obtient une micro F-mesure de 0.942 sur le jeu de test privé. Ces résultats montrent que des modèles simples tels que CRF linéaire entraînés sur des caractéristiques simples telles que la forme des mots sont efficaces pour la tâche de reconnaissance d'entités spatiales nommées et non nommées.

1 Introduction

Le défi TextMine 2024 consiste en la reconnaissance d'entités spatiales nommées et non nommées à plusieurs niveaux dans les Instructions nautiques, une série d'ouvrages publiée par le Shom (Service hydrographique et océanographique de la Marine). Le corpus mis à disposition est constitué d'extraits de 15 volumes des Instructions nautiques décrivant des environnements maritimes côtiers. La tâche de reconnaissance des entités nommées (REN) a suscité un vif intérêt ces dernières années, et a vu l'application de nombreuses approches, en particulier les approches par règles (Nadeau et Sekine, 2007; Sekine et Nobata, 2004; Kim et Woodland, 2000; Farmakiotou et al., 2000; Appelt et al., 1993), par apprentissage automatique (*machine-learning*) (Asahara et Matsumoto, 2003; Borthwick et al., 1998; Sekine, 1998; Bikel et al., 1998; McCallum et Li, 2003; Ritter et al., 2011; Rocktäschel et al., 2012), et plus récemment par apprentissage profond (*deep-learning*) (Yadav et Bethard, 2019; Li et al., 2020)

Dans cet article, nous présentons notre participation au défi TextMine 2024. Nous abordons la reconnaissance d'entités spatiales nommées et non nommées comme une tâche d'annotation de séquence. Nous entraînons pour cela trois modèles CRF linéaires (Lafferty et al., 2001), chacun entraîné sur des ensembles de caractéristiques différents. De plus, nous proposons un modèle hybride, combinant un modèle CRF linéaire pour la détection des mentions d'entités nommées, et un modèle XGBoost (Chen et Guestrin, 2016) pour la classification des mentions détectées. Nous entraînons les modèles sur des caractéristiques (*features*) morphologiques et

syntaxiques extraites de chaque mot dans le jeu de données d’entraînement mis à disposition. Les résultats du défi sont disponibles sur Kaggle¹, tandis que notre code et nos expérimentations sont disponibles sur GitHub². Le reste de l’article est structuré comme suit : nous présentons les données dans la section 2, puis notre méthodologie dans la section 3. Nous présentons l’évaluation de notre méthodologie en section 4 avant de présenter notre conclusion en section 5.

2 Données

Un jeu de données d’entraînement et de test ont été mis à disposition des participants, au format CSV. Le jeu d’entraînement est divisé en deux colonnes, une pour les tokens et une pour les classes correspondantes. Le jeu de test ne contient que les tokens. Le jeu de test mis à disposition, dit jeu de test public, représente environ 60 % des données totales de test. Afin d’entraîner nos modèles et identifier les meilleures combinaisons d’hyper-paramètres possibles, nous avons divisé le jeu d’entraînement en un jeu d’entraînement et un jeu de développement. Le Tableau 1 présente la distribution des classes dans les jeux de données d’entraînement et de développement.

Classe	Entraînement	Développement
aucun	22 859	9 804
geogFeat	1 899	805
geogName name	1 491	626
geogFeat geogName	995	468
geogName	655	255
Total	27 899	11 958

TAB. 1 – *Distribution des classes dans les jeux de données d’entraînement et de développement.*

Plusieurs classes sont composées de multiples étiquettes : par exemple, la classe *geogFeat geogName* est composée des étiquettes *geogFeat* et *geogName*, qui sont toutes deux des classes à part entière dans le jeu de données. Le Tableau 2 présente la distribution des classes uniques dans les jeux d’entraînement et de développement.

Classe	Entraînement	Développement
name	1 491	626
geogFeat	2 894	1 273
geogName	3 141	1 349

TAB. 2 – *Distribution des étiquettes uniques dans les jeux de données d’entraînement et de développement.*

1. <https://www.kaggle.com/competitions/defi-textmine-2024/submissions>

2. <https://github.com/nicolasgutehrl/DefiTextmineCRIT>

3 Méthodologie et implémentation

Nous abordons ce défi comme une tâche d'annotation de séquence. Afin d'entraîner les modèles, nous extrayons des caractéristiques morphologiques et syntaxiques pour chaque token t dans le jeu de données. Nous extrayons ces mêmes caractéristiques des token qui précèdent et suivent directement le token t . Ces caractéristiques sont préfixées par *prev_* et *next_* respectivement. Le Tableau 3 présente l'ensemble des caractéristiques extraites pour un token t . La caractéristique *shape* représente les lettres majuscules par un "X", les valeurs numériques par un "d" et les ponctuations par un ".". Tous les autres caractères sont représentés par un "x". Par exemple, l'entité nommée "Nadji" aura la forme "Xxxxx". Nous employons le modèle *fr_core_news_lg* du framework *spaCy*³ pour obtenir les parties du discours ainsi que les rôles syntaxiques des tokens. Nous employons également la liste de mots vides de ce modèle pour déterminer si un token est un mot vide.

	Caractéristique	Description
1	<i>Token</i>	Mot
2	<i>lower</i>	Mot en minuscule
3	<i>isdigit</i>	Vrai si le mot est une valeur numérique, sinon Faux
4	<i>isupper</i>	Vrai si le mot est écrit en lettres capitales, sinon Faux
5	<i>ispunct</i>	Vrai si le mot est une ponctuation, sinon Faux
6	<i>isstop</i>	Vrai si le mot est un mot vide, sinon Faux
7	<i>len</i>	Nombre de caractères composant le mot
8	<i>shape</i>	Forme du mot
9	<i>pos</i>	Partie du discours du mot
10	<i>dep</i>	Rôle syntaxique du mot

TAB. 3 – Caractéristiques extraites pour chaque token.

Nous entraînons trois modèles CRF linéaires, chacun avec un ensemble différent de caractéristiques : *baseModel* entraîné avec les caractéristiques 1 à 8, *posModel*, entraîné avec les caractéristiques 1 à 9 et *posDepModel*, entraîné avec les caractéristiques 1 à 10. Notre objectif est de comparer l'apport des parties du discours et des dépendances syntaxiques pour la reconnaissance d'entités spatiales nommées et non nommées. Nous utilisons l'implémentation du modèle CRF linéaire proposée par le package *sklearn-crfsuite*⁴. Pour chaque modèle, nous recherchons la combinaison optimale de valeurs d'hyper-paramètres parmi les valeurs présentées dans le Tableau 4. Nous évaluons chaque modèle sur le jeu de développement. Les valeurs optimales de chaque hyper-paramètre pour chaque modèle sont présentées dans le Tableau 5. Enfin, nous entraînons chaque modèle CRF linéaire sur l'ensemble des données, c'est-à-dire sur les jeux d'entraînement et de développement combinés, avec les meilleures valeurs d'hyper-paramètres identifiées.

Comme indiqué en section 2, la classe d'un token peut être une combinaison de plusieurs étiquettes. Afin d'exploiter la nature composite des classes, nous proposons un modèle hy-

3. <https://spacy.io/>

4. <https://sklearn-crfsuite.readthedocs.io/en/latest/>

c1	0.5, 0.1, 0.01, 0.001, 0.0001
c2	0.5, 0.1, 0.01, 0.001, 0.0001

TAB. 4 – Valeurs testées pour les hyper-paramètres *c1* et *c2* du modèle CRF linéaire.

Modèle	c1	c2
<i>baseModel</i>	0.1	0.001
<i>posModel</i>	0.1	0.01
<i>posDepModel</i>	0.1	0.01

TAB. 5 – Valeurs optimales des hyper-paramètres pour les modèles *baseModel*, *posModel* et *posDepModel*.

bride, qui combine un modèle CRF linéaire pour la détection d'entités nommées et un modèle XGBoost pour la classification de ces entités.

Afin d'entraîner le modèle CRF linéaire de ce modèle hybride, nous modifions les jeux d'entraînement et de développement de telle sorte que toutes les classes autre que "aucun" sont remplacées par l'étiquette "NER". Nous proposons ainsi une version binaire de jeu de données. La distribution des classes dans les jeux d'entraînement et de développement binaires est présentée dans le Tableau 6.

Classe	Entraînement	Développement
aucun	22 859	9 804
NER	5 040	2 154
Total	27 899	11 958

TAB. 6 – Distribution des classes binaires dans les jeux de données d'entraînement et de développement.

Nous entraînons plusieurs modèles CRF linéaires afin de trouver la combinaison optimale de valeurs d'hyper-paramètres parmi les valeurs présentées dans le Tableau 7. La valeur optimale de chaque hyper-paramètre y est notée en gras. Chaque modèle est évalué sur le jeu de développement. Nous entraînons enfin le modèle CRF linéaire binaire final sur l'ensemble des données avec les meilleurs valeurs d'hyper-paramètres identifiées.

Pour entraîner le modèle XGBoost, nous modifions les jeux de données d'entraînement et de développement de telle sorte qu'un token est associé à plusieurs classes. Ainsi, nous entraînons ce modèle selon une tâche classification multi-étiquettes (*multi-label*), et non pas une tâche de classification multi-classes. Nous employons l'implémentation du modèle XGBoost du package *xgboost*. Nous entraînons plusieurs modèles XGBoost afin de trouver la combinaison optimale de valeurs d'hyper-paramètres parmi les valeurs présentées dans le Tableau 8. La valeur optimale de chaque hyper-paramètre y est notée en gras. Chaque modèle est évalué sur le jeu de développement. Nous entraînons enfin le modèle XGBoost final sur l'ensemble des données avec les meilleurs valeurs d'hyper-paramètres identifiées.

c1	0.5, 0.1, 0.01, 0.001, 0.0001
c2	0.5, 0.1 , 0.01, 0.001, 0.0001

TAB. 7 – Valeurs testées pour les hyper-paramètres *c1* et *c2* du modèle CRF linéaire. Les valeurs en gras sont les valeurs optimales identifiées.

n_estimators	50, 100, 150 , 200, 350, 500
max_depth	1, 2, 3, 4, 5, 6, 7, 8, 9 , 10
learning_rate	0.0001, 0.001, 0.01, 0.1, 1.0
subsample	.25, .5, .75, 1

TAB. 8 – Valeurs testées pour les hyper-paramètres *n_estimators*, *max_depth*, *learning_rate* et *subsample* du modèle XGBoost. Les valeurs en gras sont les valeurs optimales identifiées.

4 Résultats et Discussion

Les modèles ont été évalués en terme de micro F-mesure sur un jeu de test public et un jeu de test privé. Le jeu de test public représente environ 60 % du jeu de test privé. Le Tableau 9 présente les performances des quatre modèles sur les jeux de test public et privé.

Modèle	Micro F-mesure (public)	Micro F-mesure (privé)
<i>baseModel</i>	0.964	0.942
<i>posModel</i>	0.959	0.944
<i>posDepModel</i>	0.959	0.943
<i>hybrid</i>	0.915	0.895

TAB. 9 – Evaluation des quatre modèles selon la micro F-mesure sur les jeux de test public et privé.

Sur le jeu de test public, le modèle *baseModel* atteint le meilleur score, qui est de 0.964. Les modèles *posModel* et *posDepModel* obtiennent tous deux un score de 0.959. Enfin, le modèle *hybrid* obtient le score le plus bas, qui est de 0.915. Sur le jeu de test privé, le modèle *posModel* atteint le meilleur score, qui est de 0.944. Les modèles *baseModel* et *posDepModel* obtiennent respectivement un score de 0.942 et 0.943. Enfin, le modèle *hybrid* obtient le score le plus bas, qui est de 0.895.

Les scores des modèles *baseModel*, *posModel* et *posDepModel* sont similaires. L'ajout des parties du discours et des rôles syntaxiques comme caractéristiques d'entraînement ne semblent pas affecter les performances du modèle. De simples caractéristiques d'entraînements, telles que celles que nous avons extraites, semblent donc adéquates pour identifier les entités spatiales nommées et non nommées. Cependant, les données mises à disposition rassemblent les tokens dans un seul grand document, sans les séparer en phrases. Cela peut affecter la qualité de l'analyse en partie du discours ainsi que l'analyse en dépendances syntaxiques. Le modèle *hybrid* obtient les scores les plus bas sur les deux évaluations. Cette baisse des scores peut s'expliquer par l'incapacité du modèle XGBoost à prendre en compte le

contexte de l'élément à catégoriser. La conversion de la tâche de classification multiclassées en une tâche de classification multi-étiquettes peut également expliquer ces résultats inférieurs.

Le Tableau 10 présente les dix plus importantes caractéristiques identifiées par le modèle *baseModel* contribuant positivement à l'identification de chaque classe. Chaque caractéristique est représentée par sa valeur et son poids. La caractéristique *lower*, c'est-à-dire la forme du mot en minuscule, est la plus fréquemment employée pour chaque classe. Le modèle se concentre sur le mot lui-même pour les classes contenant l'étiquette *geogFeat*, comme en témoigne l'emploi des caractéristiques *Token* et *lower*. Il en est de même pour les classes contenant l'étiquette *geogName* et *name*. Pour ces classes là, le modèle se concentre également sur les mots précédents et suivants, comme en témoigne l'emploi des caractéristiques *prev-Token*, *prev_lower*, *next-Token*, *next_lower*. Le modèle semble apprendre du vocabulaire maritime ("mouillage", "rade", "lagon", etc.) pour identifier les classe *geogFeat* et *geogFeat geogName*. De même, il semble apprendre des noms propres désignant des lieux ("Islande", "Djibouti", etc.) pour identifier les classes *geogFeat geogName* et *geogName name*.

aucun	geogName	geogName name	geogFeat	geogFeat geogName
<i>isdigit</i> , 8.512	<i>Token</i> :», 6.085	<i>Token</i> :nouveau, 7.609	<i>lower</i> :mouillages, 11.987	<i>lower</i> :île, 10.255
<i>lower</i> :entre, 7.967	<i>lower</i> :», 6.085	<i>lower</i> :nouveau, 7.609	<i>lower</i> :amers, 9.929	<i>lower</i> :pointe, 8.422
<i>Token</i> :petit, 7.822	<i>Token</i> :principal, 4.995	<i>prev_lower</i> :côte, 7.492	<i>lower</i> :rade, 9.815	<i>lower</i> :ras, 7.632
<i>lower</i> :ouest, 7.780	<i>lower</i> :principal, 4.995	<i>lower</i> :djibouti, 7.405	<i>lower</i> :lagon, 9.778	<i>lower</i> :cap, 6.950
<i>Token</i> :grande, 6.443	<i>prev_lower</i> :capitainerie, 4.939	<i>lower</i> :blancpignon, 7.383	<i>lower</i> :alignement, 9.282	<i>lower</i> :îles, 6.755
<i>lower</i> :petite, 6.408	<i>next-Token</i> :médián, 4.922	<i>lower</i> :grande, 7.100	<i>Token</i> :appontement 9.074	<i>lower</i> :aéroport, 6.554
<i>lower</i> :jusqu', 5.881	<i>next_lower</i> :médián, 4.922	<i>lower</i> :espagnole, 6.454	<i>lower</i> :pontons, 9.057	<i>Token</i> :Island, 6.535
<i>shape</i> :X, 5.880	<i>shape</i> :xxxxxxx, 4.187	<i>Token</i> :ancien, 6.296	<i>Token</i> :baie, 8.515	<i>lower</i> :island, 6.535
<i>lower</i> :devant, 5.782	<i>prev_lower</i> :phare, 4.123	<i>lower</i> :ancien, 6.296	<i>lower</i> :échelles, 8.266	<i>Token</i> :Pass, 6.496
<i>Token</i> :Est, 5.746	<i>next-Token</i> :sur, 3.932	<i>shape</i> :XXXXX.Xxxxxx, 6.258	<i>lower</i> :îles, 8.249	<i>prev-Token</i> :ancienne, 6.251

TAB. 10 – Dix plus importantes caractéristiques identifiées par le modèle *baseModel* contribuant positivement à l'identification de chaque classe. Chaque caractéristique est représentée par sa valeur et son poids.

5 Conclusion

Dans cet article, nous avons présenté notre participation au Défi TextMine 2024 : "Reconnaissance d'entités géographiques dans un corpus des Instructions nautiques". Nous avons proposé deux méthodes pour la tâche de reconnaissance d'entités spatiales nommées et non nommées. Notre première méthode repose sur un modèle CRF linéaire, tandis que notre seconde méthode est hybride, et combine un modèle CRF linéaire avec un modèle XGBoost. Notre modèle CRF linéaire, entraîné sur des caractéristiques morphologiques simples telles que la forme des mots, obtient une micro F-mesure de 0.964 sur le jeu de test public, et une micro F-mesure de 0.942 sur le jeu de test privé. Ces résultats montrent que des modèles simples tels que CRF linéaire entraînés sur des caractéristiques morphologiques également simples sont efficaces pour la tâche de reconnaissance d'entités spatiales nommées et non nommées. Dans des travaux futurs, nous comptons expérimenter avec d'autres caractéristiques d'entraînement d'ordre morphologique, syntaxique ou sémantique. Nous avons également l'intention d'approfondir notre étude sur la nature composite des entités nommées et non-nommées en continuant nos expérimentations sur la classification multi-étiquettes.

Références

- Appelt, D. E., J. R. Hobbs, J. Bear, D. Israel, et M. Tyson (1993). Fastus : A finite-state processor for information extraction from real-world text. In *IJCAI*, Volume 93, pp. 1172–1178.
- Asahara, M. et Y. Matsumoto (2003). Japanese named entity extraction with redundant morphological analysis. In *Proceedings of the 2003 human language technology conference of the North American chapter of the association for computational linguistics*, pp. 8–15.
- Bikel, D. M., S. Miller, R. Schwartz, et R. Weischedel (1998). Nymble : a high-performance learning name-finder. *arXiv preprint cmp-lg/9803003*.
- Borthwick, A., J. Sterling, E. Agichtein, et R. Grishman (1998). Description of the mene named entity system as used in muc-7. In *Proceedings of the Seventh Message Understanding Conference (MUC-7), Fairfax, Virginia, April 29-May 1, 1998*.
- Chen, T. et C. Guestrin (2016). Xgboost : A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794.
- Farmakiotou, D., V. Karkaletsis, J. Koutsias, G. Sigletos, C. D. Spyropoulos, et P. Stamatopoulos (2000). Rule-based named entity recognition for greek financial texts. In *Proceedings of the Workshop on Computational lexicography and Multimedia Dictionaries (COMLEX 2000)*, pp. 75–78.
- Kim, J.-H. et P. C. Woodland (2000). A rule-based named entity recognition system for speech input. In *Sixth International Conference on Spoken Language Processing*.
- Lafferty, J., A. McCallum, et F. C. Pereira (2001). Conditional random fields : Probabilistic models for segmenting and labeling sequence data.
- Li, J., A. Sun, J. Han, et C. Li (2020). A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering* 34(1), 50–70.
- McCallum, A. et W. Li (2003). Early results for named entity recognition with conditional random fields, feature induction and web-enhanced lexicons.
- Nadeau, D. et S. Sekine (2007). A survey of named entity recognition and classification. *Linguisticae Investigationes* 30(1), 3–26.
- Ritter, A., S. Clark, O. Etzioni, et al. (2011). Named entity recognition in tweets : an experimental study. In *Proceedings of the 2011 conference on empirical methods in natural language processing*, pp. 1524–1534.
- Rocktäschel, T., M. Weidlich, et U. Leser (2012). Chemspot : a hybrid system for chemical named entity recognition. *Bioinformatics* 28(12), 1633–1640.
- Sekine, S. (1998). Nyu : Description of the japanese ne system used for met-2. <http://www.muc.saic.com/>.
- Sekine, S. et C. Nobata (2004). Definition, dictionaries and tagger for extended named entity hierarchy. In *LREC*, pp. 1977–1980. Lisbon, Portugal.
- Yadav, V. et S. Bethard (2019). A survey on recent advances in named entity recognition from deep learning models. *arXiv preprint arXiv :1910.11470*.

Summary

This article presents our participation to the TextMine 2024 challenge. We propose a method based on a linear CRF model, as well as a hybrid method combining a linear CRF model and an XGBoost model. We train the models on morphological and syntactic features extracted from each word. Our linear CRF model obtains a micro F-measure of 0.942 on the private test set. These results show that simple models such as linear CRF trained on simple features such as word shape are effective for the Named and Unnamed Spatial Entity Recognition task.