

Identifying professions & occupations in Health-related Social Media using Natural Language Processing

J. Alberto Mesa Murgado and Ana Belén Parras Portillo and Pilar López-Úbeda
and M. Teresa Martín-Valdivia and L. Alfonso Ureña López

SINAI Research Group - CEATIC - Universidad de Jaén

Campus Las Lagunillas s/n. E-23071, Jaén, Spain

{jmurgado, abparras, plubeda, maite, laurena}@ujaen.es

Abstract

This paper describes the entry of the research group SINAI at SMM4H's ProfNER task on identifying professions and occupations in social media data related to health. Specifically, we participated in Task 7a: Tweet Binary Classification to determine whether a tweet contains mentions of occupations or not and also in Task 7b: NER Offset Detection and Classification aimed at predicting occupations mentions and classify them as either professions or working statuses.

1 Introduction

Natural Language Processing (NLP) and Machine Learning (ML) techniques are becoming essential in critical fields such as the one of healthcare, considering that they perform tasks faster than a human agent and at a very high level of reliability. Some of these tasks include the automatic assignment of International Classification of Diseases (ICD) codes to health related texts (Perea-Ortega et al., 2020) or the detection of negative and positive emotions in medical documents (Plaza-del Arco et al., 2019).

Automatic text classification and Named Entity Recognition (NER) are two tasks in which NLP has proved to have a relevant impact. In both cases we are given a certain set of documents and while for the first task we aim to classify them distinguishing by a certain criteria, the second seeks to detect and tag specific entities.

The Social Media Mining for Health (SMM4H) 2021 ProfNER Shared Task (Miranda-Escalada et al., 2021) emphasizes the importance of identifying professions and occupations within social media content related to health, this knowledge could later be applied to determine which of them are at risk due to direct exposure to the COVID-19 pandemic and/or state what professional sectors are more prone to mental health issues due to the uncertainty of the current situation.

This issue has been further subdivided into two tracks: determine whether the social media textual content contains mentions of professions or not (a binary classification task) and to identify professions and working statuses within the text in order to extract its text span and tag it accordingly (NER). Our research group has used NLP and ML approaches for both tasks in combination with two dictionaries which we have developed.

Considering this information, this paper is structured as follows: section 2 introduces work related to this challenge and research field. Section 3 briefly describes the dataset provided and its characteristics. Section 4 states the systems we have developed for each task. Section 5 exhibits the results from our systems using the test dataset and finally, in Section 6 we present our conclusions and future work.

2 Related work

Social media plays an important role where people can share information related to health. This information can be used for public health monitoring tasks through the use of NLP techniques.

On one hand, in terms of document classification in the medical field, many researchers have used social networks as a source of information to develop and evaluate systems. For example, to predict mental illnesses such as depression or anorexia (Al-darwish and Ahmad, 2017; López-Úbeda et al., 2021) and to detect nonmedical prescription medication (Al-Garadi et al., 2021). More recently, new studies analyzed health, psychosocial, and social issues emanating from the COVID-19 pandemic from social network comments using NLP (López-Úbeda et al., 2020a; Müller et al., 2020; Oyebode et al., 2020).

On the other hand, several NER systems have been developed using NLP-based systems such as MedLEE (Friedman, 1997), MetaMap (Aronson and Lang, 2010) and cTAKES (Savova et al.,

2010)). Most of these are rule-based systems that use extensive medical vocabularies. Current state-of-the-art approaches to the NER task propose the use of RNNs to learn useful representations automatically because they facilitate the modeling of long-distance dependencies between words in a sentence (López-Úbeda et al., 2019; López-Úbeda et al., 2020b).

Since there is currently a great growth in demand for classification and extraction of information from medical texts, the NLP community has organized a series of open challenges with a focus on biomedical entity extraction and document classification tasks such as DDIExtraction (Segura Bedmar et al., 2013), the N2C2 - National NLP Clinical Challenges shared task (Henry et al., 2020) and the CHEMDNER challenge (Krallinger et al., 2015). Finally, SMM4H (Weissenbacher et al., 2019) provided tasks for the extraction of adverse effects using Twitter as a source of information. In this workshop, participants were first required to identify whether a tweet contained an Adverse Drug Reaction (ADR). Subsequently, the challenge provided the NER task to locate the specific ADR.

The use of Spanish as the main language of a challenge has emerged in recent years providing important workshops such as the DIANN (Fabregat et al., 2018) (Disability Annotation Task) task, PharmaCoNER (Agirre et al., 2019) (Pharmacological Substances, Compounds and proteins and NER), Cantemist (Miranda-Escalada et al., 2020) and eHealth-KD (Piad-Morffis et al., 2020) (eHealth knowledge discovery).

3 Dataset

Organizers provided us with a dataset consisting of 8,000 tweets from Twitter subdivided into two subsets: 6,000 tweets with which to train our systems and 2,000 tweets to validate them. Namely, train set and dev set, accordingly.

Besides the tweet's identifier and text span, a binary value was used to determine whether it contained a profession or not as well as its corresponding annotated entities tagged using the Inside-outside-beginning (IBO) format.

To compare the performance of all the systems presented at this shared task, we were provided with another dataset, namely test set, consisting of 2,000 processed tweets and 25,000 raw tweets for the background set.

4 Methodology

For our participation we employ ML models enriched with custom made dictionaries consisting of 776 professions such as "Farmacéutico" (Pharmacist), "Dentista" (Dentist), "Cajera" (Cashier) and "Veterinario" (Vet) recovered from the "Listado de profesiones reguladas en el ámbito sanitario"¹ provided by the Ministerio de de Sanidad y Política Social of the Spanish Government and from the occupations listed by the European Commission in their International Standard Classification of Occupations (ISCO)².

The second dictionary contains 26 working statuses including "Autónomo", "Funcionario" (Public employee) and "Erte" (record of temporary Employment regulation) among others based on the Workshop's annotation guidelines³.

4.1 Task 7a: Binary Classification

To classify tweets whether if they contained mentions of professions or not, we used two approaches: a Support Vector Machine (SVM) and bidirectional Long Short Term Memory (BiLSTM) Recurrent Neural Network (RNN), both combined with our professions dictionary.

The SVM approach applies the scikit-learn library (Pedregosa et al., 2011) using its default parameters, as stated in the documentation, the words for each tweet, also stated as document, were transformed into vectors considering the frequency of the terms within each document (TFIDF). This structure was later enriched using our professions dictionary through a vector for each document consisting of as much binary values as the number of terms in the dictionary such as each binary value represented if the term within the document was in the dictionary or not.

The BiLSTM model is implemented using the Tensorflow library (Abadi et al., 2015) which we explored through different batch sizes, ranging from 2^7 (2^6 and 2^8), and different number of epochs. Although the best accuracy obtained for training (0.8695) determined a batch size of 128 and 5 epochs. This model makes use of the GloVe (Pennington et al., 2014) word embeddings 200d vector,

¹https://www.msrebs.gob.es/eu/profesionales/formacion/docs/Anexo_X_del_Real_Decreto_1837.pdf

²<https://ec.europa.eu/esco/portal/occupation>

³<https://zenodo.org/record/4306017#.YE3vpJ1KhhE>

pre-trained using Twitter’s tweets, this approach is also uses our professions dictionary.

4.2 Task 7b: NER offset detection and classification

To detect and classify entities within those same tweets and retrieve them discerning the text span of the entities as well as their initial and end position within the text, we opted for a Conditional Random Fields (CRF) approach.

This approach was implemented using scikit-learn crfsuite library (Wijffels and Okazaki, 2007-2018), transforming words to features and using the words that came before and after each term, then enriched the system using both our dictionaries (added features as binary values). For this implementation we considered the L-BFGS method for the gradient descent, a 100 iterations and values 0.1 for both c1 and c2.

Our system assigned an IBO tag to each term of every tweet and we later searched for the entity text span within the same tweet to extract its initial and end position.

5 Evaluation results

We have been provided with the median of the participants’ score using three metrics: precision (P), recall (R) and F1-scoring (F1) (Magee et al., 2021). Using them we built 2 tables consisting of two sections: the upper section shows the score obtained by our systems on the test set, while the latter fits the same purpose for the dev set.

5.1 Task 7a: Binary Classification

We submitted 2 runs for the evaluation phase: a combination of SVM (SVM+Dic) and an BiLSTM model (BiLSTM+Dic), both combined with our professions dictionary, the scoring associated to these systems applied to the provided datasets is displayed in Table 1.

Model	P	R	F1
Median	0.9185	0.8553	0.85
BiLSTM+Dic	0.7612	0.4752	0.59
SVM+Dic	0.8995	0.4255	0.58
BiLSTM+Dic	0.77	0.72	0.74
SVM+Dic	0.86	0.70	0.74
SVM	0.86	0.64	0.66

Table 1: Scores obtained by our systems on the SMM4H ProfNER Shared Task - Task 7a (binary classification) applied over the test and dev set, accordingly.

While the performance of our systems on the training dataset was, at average, close to 0.74 (F1), on the test set it was decreased in a 21%. Therefore, resulting in a value close to 0,59 (again, F1), 31% below the median.

5.2 Task 7b: NER Offset Classification

We were closer to the median and in line with what our systems obtained on the training dataset (1% decrease in performance compared to the test set). These results are displayed in Table 2 in the same way in which Table 1 was: the upper section refers to the scoring for our systems on the test set while the latter, exhibits the scoring for the dev set.

Model	P	R	F1
Median	0.842	0.7265	0.7605
CRF+Dic	0.824	0.652	0.728
CRF+Dic	0.861	0.647	0.739
CRF	0.852	0.597	0.702

Table 2: Score obtained by our systems on the SMM4H ProfNER Shared Task - Task 7b (NER) applied over the test and dev set, accordingly.

6 Conclusion

For our participation in the SMM4H Task 7 ProfNER Shared Task on identifying professions and occupations we implemented three systems. First two are aimed at the binary classification task (Task A) using an SVM and a BiLSTM approach, both combined with our professions dictionary. The latter system follows a CRF approach combined with our professions and working statuses dictionaries and it is applied to the NER task (Task B).

Our predictions for the training set were consistent with those obtained on the test set for the second task (NER) whereas our approaches for the first task (binary classification) fell short of our expectations by 21% below our training results.

For future work, we will use the gold test in order to perform a deeper analysis to assess why did this event happened and therefore improve the performance of our systems.

Acknowledgements

This work has been partially supported by the LIVING-LANG project [RTI2018-094653-B-C21] of the Spanish Government and the Fondo Europeo de Desarrollo Regional (FEDER).

References

- 248
- 249 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene
250 Brevdo, Zhifeng Chen, Craig Citro, Greg S. Cor-
251 rado, Andy Davis, Jeffrey Dean, Matthieu Devin,
252 Sanjay Ghemawat, Ian Goodfellow, Andrew Harp,
253 Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal
254 Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh
255 Levenberg, Dan Mané, Rajat Monga, Sherry Moore,
256 Derek Murray, Chris Olah, Mike Schuster, Jonathon
257 Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar,
258 Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan,
259 Fernanda Viégas, Oriol Vinyals, Pete Warden, Mar-
260 tin Wattenberg, Martin Wicke, Yuan Yu, and Xiao-
261 qiang Zheng. 2015. [TensorFlow: Large-scale ma-
262 chine learning on heterogeneous systems](#). Software
263 available from tensorflow.org.
- 264 Aitor Gonzalez Agirre, Montserrat Marimon, Ander In-
265 txaurrondo, Obdulia Rabal, Marta Villegas, and Mar-
266 tin Krallinger. 2019. Pharmaconer: Pharmacologi-
267 cal substances, compounds and proteins named en-
268 tity recognition track. In *Proceedings of The 5th
269 Workshop on BioNLP Open Shared Tasks*, pages 1–
270 10.
- 271 Mohammed Ali Al-Garadi, Yuan-Chi Yang, Haitao
272 Cai, Yucheng Ruan, Karen O’Connor, Gonzalez-
273 Hernandez Graciela, Jeanmarie Perrone, and Abeer
274 Sarker. 2021. Text classification models for the auto-
275 matic detection of nonmedical prescription medica-
276 tion use from social media. *BMC medical informat-
277 ics and decision making*, 21(1):1–13.
- 278 M. M. Aldarwish and H. F. Ahmad. 2017. [Predicting
279 depression levels using social media posts](#). In *2017
280 IEEE 13th International Symposium on Autonomous
281 Decentralized System (ISADS)*, pages 277–280.
- 282 Alan R Aronson and François-Michel Lang. 2010. An
283 overview of metamap: historical perspective and re-
284 cent advances. *Journal of the American Medical In-
285 formatics Association*, 17(3):229–236.
- 286 Hermenegildo Fabregat, Juan Martinez-Romo, and
287 Lourdes Araujo. 2018. Overview of the DIANN
288 Task: Disability Annotation Task. In *IberEval@ SE-
289 PLN*, pages 1–14.
- 290 Carol Friedman. 1997. Towards a comprehensive medi-
291 cal language processing system: methods and issues.
292 In *Proceedings of the AMIA annual fall symposium*,
293 page 595. American Medical Informatics Associa-
294 tion.
- 295 Sam Henry, Kevin Buchan, Michele Filannino, Amber
296 Stubbs, and Ozlem Uzuner. 2020. 2018 n2c2 shared
297 task on adverse drug events and medication extrac-
298 tion in electronic health records. *Journal of the
299 American Medical Informatics Association*, 27(1):3–
300 12.
- 301 Martin Krallinger, Florian Leitner, Obdulia Rabal,
302 Miguel Vazquez, Julen Oyarzabal, and Alfonso Va-
303 lencia. 2015. Chemdner: The drugs and chemical
names extraction challenge. *Journal of cheminfor-
matics*, 7(1):S1.
- Pilar López-Úbeda, Manuel Carlos Díaz-Galiano,
Teodoro Martín-Noguerol, Antonio Luna, L Alfonso
Ureña-López, and M Teresa Martín-Valdivia. 2020a.
Covid-19 detection in radiological text reports inte-
grating entity recognition. *Computers in Biology
and Medicine*, 127:104066.
- Pilar López-Úbeda, Manuel Carlos Díaz Galiano,
M Teresa Martín-Valdivia, and L Alfonso Urena
Lopez. 2019. Using machine learning and deep
learning methods to find mentions of adverse drug
reactions in social media. In *Proceedings of the
Fourth Social Media Mining for Health Applications
(# SMM4H) Workshop & Shared Task*, pages 102–
106.
- Pilar López-Úbeda, José M Perea-Ortega, Manuel C
Díaz-Galiano, M Teresa Martín-Valdivia, and L Al-
fonso Ureña-López. 2020b. Sinai at ehealth-kd
challenge 2020: Combining word embeddings for
named entity recognition in spanish medical records.
- Pilar López-Úbeda, Flor Miriam Plaza-del Arco,
Manuel Carlos Díaz-Galiano, and Maria-Teresa
Martín-Valdivia. 2021. How successful is transfer
learning for detecting anorexia on social media? *Ap-
plied Sciences*, 11(4):1838.
- Arjun Magge, Ari Klein, Ivan Flores, Ileyar Al-
imova, Mohammed Ali Al-garadi, Antonio Miranda-
Escalada, Zulfat Miftahutdinov, Eulàlia Farré-
Maduell, Salvador Lima López, Juan M Banda,
Karen O’Connor, Abeer Sarker, Elena Tutubalina,
Martin Krallinger, Davy Weissenbacher, and Gra-
ciela Gonzalez-Hernandez. 2021. Overview of the
sixth social media mining for health applications (#
smm4h) shared tasks at naacl 2021. In *Proceedings
of the Sixth Social Media Mining for Health Appli-
cations Workshop & Shared Task*.
- A Miranda-Escalada, E Farré, and M Krallinger. 2020.
Named entity recognition, concept normalization
and clinical coding: Overview of the cantemist
track for cancer text mining in spanish, corpus,
guidelines, methods and results. In *Proceedings of
the Iberian Languages Evaluation Forum (IberLEF
2020)*, *CEUR Workshop Proceedings*.
- Antonio Miranda-Escalada, Eulàlia Farré-Maduell, Sal-
vador Lima López, Vicent Briva-Iglesias, Marvin
Agüero-Torales, Luis Gascó-Sánchez, and Martin
Krallinger. 2021. The profner shared task on
automatic recognition of professions and occupa-
tion mentions in social media: systems, evaluation,
guidelines, embeddings and corpora. In *Proceed-
ings of the Sixth Social Media Mining for Health Ap-
plications Workshop & Shared Task*.
- Martin Müller, Marcel Salathé, and Per E Kummervold.
2020. Covid-twitter-bert: A natural language pro-
cessing model to analyse covid-19 content on twitter.
arXiv preprint arXiv:2005.07503.

361 Oladapo Oyeboade, Chinenye Ndulue, Ashfaq Adib, Di- 417
362 nesh Mulchandani, Banuchitra Suruliraj, Fidelia An- 418
363 ulika Orji, Christine Chambers, Sandra Meier, and 419
364 Rita Orji. 2020. Health, psychosocial, and social is-
365 sues emanating from covid-19 pandemic based on
366 social media comments using natural language pro-
367 cessing. *arXiv preprint arXiv:2007.12144*.

368 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel,
369 B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,
370 R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,
371 D. Cournapeau, M. Brucher, M. Perrot, and E. Duch-
372 esnay. 2011. Scikit-learn: Machine learning in
373 Python. *Journal of Machine Learning Research*,
374 12:2825–2830.

375 Jeffrey Pennington, Richard Socher, and Christopher D.
376 Manning. 2014. *Glove: Global vectors for word rep-*
377 *resentation*. In *Empirical Methods in Natural Lan-*
378 *guage Processing (EMNLP)*, pages 1532–1543.

379 José M Perea-Ortega, Pilar López-Úbeda, Manuel C
380 Díaz-Galiano, M Teresa Martín-Valdivia, and L Al-
381 fonso Ureña-López. 2020. Sinai at clef ehealth 2020:
382 testing different pre-trained word embeddings for
383 clinical coding in spanish.

384 Alejandro Piad-Morffis, Yoan Gutiérrez, Hian
385 Cañizares-Diaz, Suilan Estevez-Velarde, Rafael
386 Muñoz, Andres Montoyo, Yudivian Almeida-Cruz,
387 et al. 2020. Overview of the ehealth knowledge
388 discovery challenge at iberlef 2020. CEUR.

389 Flor Miriam Plaza-del Arco, M Dolores Molina-
390 González, M Teresa Martín-Valdivia, and L Alfonso
391 Ureña-López. 2019. Sinai at semeval-2019 task
392 3: Using affective features for emotion classifica-
393 tion in textual conversations. In *Proceedings of the*
394 *13th International Workshop on Semantic Evalua-*
395 *tion*, pages 307–311.

396 Guergana K Savova, James J Masanz, Philip V Ogren,
397 Jiaping Zheng, Sunghwan Sohn, Karin C Kipper-
398 Schuler, and Christopher G Chute. 2010. Mayo clin-
399 ical text analysis and knowledge extraction system
400 (ctakes): architecture, component evaluation and ap-
401 plications. *Journal of the American Medical Infor-*
402 *matics Association*, 17(5):507–513.

403 Isabel Segura Bedmar, Paloma Martínez, and María
404 Herrero Zazo. 2013. Semeval-2013 task 9: Ex-
405 traction of drug-drug interactions from biomedical
406 texts (ddiextraction 2013). Association for Compu-
407 tational Linguistics.

408 Davy Weissenbacher, Abeed Sarker, Arjun Magge,
409 Ashlynn Daughton, Karen O’Connor, Michael Paul,
410 and Graciela Gonzalez. 2019. Overview of the
411 fourth social media mining for health (smm4h)
412 shared tasks at acl 2019. In *Proceedings of the*
413 *fourth social media mining for health applications*
414 *(#SMM4H) workshop & shared task*, pages 21–30.

415 Jan Wijnffels and Naoaki Okazaki. 2007-2018. *crfsuite:*
416 *Conditional random fields for labelling sequential*