# Using Twitter Data and Natural Language Processing to Generate Wildfire Heat Maps

João Cabral Pinto
cabralpinto@student.uc.pt

Catarina Silva
catarina@dei.uc.pt

Hugo Gonçalo Oliveira
hroliv@dei.uc.pt

Alberto Cardoso
alberto@dei.uc.pt

University of Coimbra
Department of Informatics Engineering
Centre for Informatics and Systems of the University of Coimbra (CISUC)
Coimbra, Portugal

## Abstract

This paper proposes an information extraction pipeline capable of generating a wildfire heat map of Portugal from Twitter posts written in Portuguese. It uses a fine-tuned version of the state-of-the-art language model BERTimbau (Portuguese BERT) to extract fire reports from large batches of recent fire-related tweets as well as the spaCy NLP library to query the location of each recently reported fire. Wildfire locations are plotted to a colored map indicating the most probable fire locations, which could prove useful in the process of allocating firefighting resources for those region. The system should be easily adaptable to work with any other country or language, provided compatible BERT and spaCy models exist.

## 1 Introduction

In recent years, wildfires have become increasingly frequent, large, and severe[1], resulting in great loss of human life, as well as the destruction of natural ecosystems and man-made infrastructure. It is therefore of great importance to create early detection mechanisms in order to reduce their spread and negative effects.

Traditionally, wildfire detection methods rely on physical devices capable of sensing heat or other signs indicative of a fire burst. In many cases, though, this approach may prove ineffective, as it may require significant governmental investment to cover all susceptible areas. In contrast, social media offers a large amount of cheaply accessible information provided by users all across the globe, updated every moment. For this reason, it has become increasingly used in the field of disaster detection, and, more specifically, fire burst detection.

In the architecture described in this work, NLP techniques are applied to large batches of recently posted tweets to generate a wildfire heat map of a country, which may help fire departments with the timely allocation of personnel to areas in need. Section 2 starts off by providing a quick rundown of some of the work published in the area of fire burst detection. In Section 3, each step of the proposed pipeline is detailed. Section 4 follows with an empiric evaluation of the pipeline. The article is concluded in Section 5, where the main takeaways of this work are discussed, as well as some of the possible tweaks it could benefit from.

## 2 Related Work

A significant amount of work has been done in the field of automatic wildfire location detection using social media. In this section, two of the most relevant articles for this work are briefly touched on. Unlike this work, they report architectures built for the English language, although the same principles should apply.

Boulton *et al.* [1] show that social media activity about wildfires is both temporally and spatially correlated with wildfire events in the contiguous United States. This is done by performing a statistical analysis on geotagged Twitter and Instagram posts containing fire-related words or hashtags such as "wildfire" and "forest fire" and two sources of wildfire occurrence data such as MODIS and FPA.

Thanos *et al.* [2] developed a pipeline more or less similar to the one proposed in this work. The authors use a hybrid LSTM-CNN classifier to determine whether a given fire-related tweet is valid or fake. Each tweet deemed valid is transformed into a tagged sentence form by applying a set of NLP procedures, namely, tokenization, part-of-speech tagging, entity recognition, and relation recognition. Fire locations are extracted from the resulting tagged sentences using a set of regular expressions defined in the training stage. At this point in the pipeline, some of the locations extracted may still not correspond to real fires, prompting the authors to propose aggregating the location points using DBSCAN and filtering out invalid reports by applying a mathematical reliability model.

## 3 Architecture

In this section, the proposed architecture, comprised of three main stages, is detailed. Firstly, the Twitter database is queried to obtain posts containing fire-related words. The obtained collection is further filtered with a fine-tuned BERT classifier to solely contain fire reports. Lastly, each fire report is geoparsed to obtain the location of each fire and mark it on the map for the user to inspect. This pipeline is summarized in Figure 1 and the corresponding Python code was also made available[2].

### 3.1 Data Fetching

Twitter is a microblogging and social networking platform with 396.5 million users globally sending 8,817 tweets every second[3], making it an ideal platform to perform social sensing. It offers a public Application Programming Interfaces (API) through which to query its database, though it is inadequate for large data collection tasks, since it limits the number of extracted tweets to 450 every 15 minutes. Instead, the unofficial API SNScrape[4] was preferred for this work, as it is able to extract a virtually unlimited amount of tweets for a given query by web scraping Twitter's search page. In order to obtain data containing wildfire reports, tweets were queried to be written in Portuguese and contain one or more fire related keywords, such as "fogo" and "incêndio", which are approximate Portuguese word equivalents to "fire" and "wildfire".

### 3.2 Classification

Fetching tweets containing fire-related words alone is not enough to restrict the collected data to fire reports, since no keyword is strictly used in posts reporting fires. This may be due the keyword having more than one connotation in the language – such as "fogo", which is also used as an interjection in Portuguese – or because the tweet it is in mentions a fire but is not a fire report – as in "Houve um grande incêndio nos anos 90" (Translation: "There was a great fire in the 90s").

This makes it necessary to employ some sort of model capable of distinguishing between real fire reports from these particular cases. To this end, after some preliminary testing, the Portuguese version of the previously described language model BERT emerged as the obvious choice. BERT was fine-tuned for the purpose by adding a feed-forward layer downstream of the pre-trained model and training it to perform the classification task. Compared to training an entire model from scratch, fine-tuning has the significant advantage of massively reducing training times without any significant loss in performance.

---

[1] https://www.epa.gov/climate-indicators/climate-change-indicators-wildfires (accessed 2022-09-27)

[2] https://github.com/cabralpinto/wildfire-heat-map-generator (accessed 2022-09-27)

[3] https://www.bestproxyreviews.com/twitter-statistics (accessed 2022-09-23)

[4] https://github.com/JustAnotherArchivist/snscrape (accessed 2022-09-27)
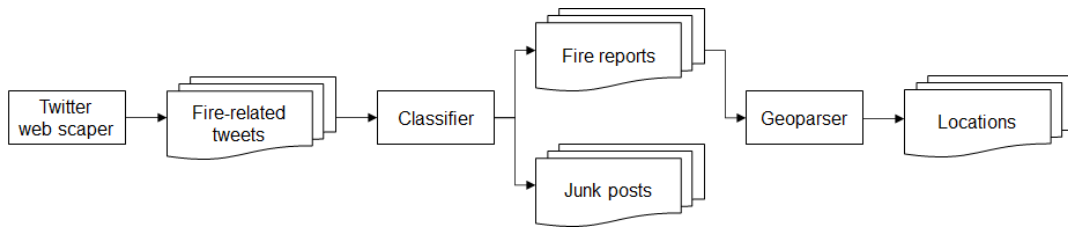
Figure 1: Proposed architecture.

## 3.3 Geoparsing

In this stage, the fire location is extracted from each fire report to then include it in the heat map. Assuming the fire report was tweeted near the location of the fire, a simple technique to achieve this would be to use the geographical metadata in the tweet in question, as is done in [1]. This method, however, poses two issues: i) a significant portion of tweets containing fire reports are from news agencies, whose locations are fixed and do not coincide with the locations of the fires they report; ii) only a very minor portion of tweets are in fact geotagged (0.51%) [1], meaning a large majority of the queried tweets, possibly containing valuable information, would not contribute to the final result. The solution is to once again resort to a language model, this time to extract the location of reported fires from the text of the corresponding tweets.

To simplify this task, the assumption is made that any location name mentioned in a fire report corresponds to the location of the fire, which seems to be quite often the case. Taking this premise as true, fire locations can be extracted through geoparsing, which is the process of parsing text to detect terms associated with geographic places. Generally, the first step of geoparsing is Named Entity Recognition (NER), which is the process of locating and classifying named entities mentioned in unstructured text into pre-defined categories, including locations. NER is applied to the text of a tweet using the spaCy Python library and results in a set of location names which are concatenated to form a preliminary geocode. The geocode is sent as a search request to the Nominatim API [5] to obtain the geometry of the corresponding region in the world map.

The final step of the pipeline is to detect intersections between the regions extracted from each fire report and the regions in a predefined area of interest such as Portugal, in the case of this work. Regions with a higher intersection count appear more orange in the resulting map, meaning there is a high volume of recent tweets reporting nearby fires. Figure 2 depicts an example of a generated heat map. Some districts in the map are depicted in brighter orange tones, which coincides with a higher number of fires being reported in these areas.
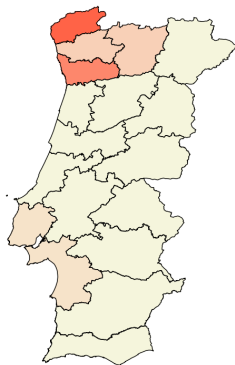


Figure 2: Heatmap generated at a time where the districts depicted in brighter oranges were suffering from a high number of wildfires.

## 4 Experimentation and Results

In this section, the classification portion of the pipeline, which uses a fine-tuned BERT model to distinguish tweets containing fire reports from other post types, is evaluated. To this end, it is compared with a baseline SVM classifier trained on the same data.

Due to the lack of a better-suited data set, all training and testing was done on a slightly modified version of a data set previously created in the

context of the FireLoc project[6] containing news headlines from various Portuguese media. In this version, each of the 2263 entries is annotated with a label regarding whether it is a fire report or not, as well as a fire location, if applicable.

Each model was trained and tested 30 times on the same random 90%/10% data set splits to obtain the unweighted mean scores shown in Table 1. As expected, fine-tuned BERT is the superior model as it achieves similar precision but much higher recall than the baseline classifier. The BERT instance with the highest F1-score (94.5%) was chosen to be used in the pipeline.

Table 1: Comparison between results of SVM and BERT.

| Classifier | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| SVM | $93.2 \pm 0.3$ | $84.9 \pm 1.4$ | $46.4 \pm 1.7$ | $59.4 \pm 1.5$ |
| BERT | $96.3 \pm 0.2$ | $83.0 \pm 1.5$ | $81.9 \pm 1.8$ | $82.1 \pm 1.4$ |

## 5 Conclusions and Future Work

In this work, an NLP pipeline was developed to transform batches of recent fire-related Tweets into wildfire heat maps. It could aid fire departments by quickly detecting fire bursts and allowing faster action.

The classifier component of the pipeline distinguishes tweets containing fire reports from other tweet categories. The state-of-the-art language model BERT was fine-tuned for this task and achieves a much higher recall than an SVM baseline model trained on the same data, suggesting that it is the preferable option in NLP classification tasks, especially with imbalanced data.

Given the opportunity to further develop this work, it is a must to develop a comprehensive data set containing tweets containing fire-related words instead of news headlines. Such data set would better suit the data queried from Twitter and should therefore lead to higher classification scores. Additionally, each reported fire in the data set should be annotated with a geometric shape representing the affected area, allowing us to empirically evaluate the geoparsing stage of the pipeline.

## Acknowledgements

## References

[1] Chris A. Boulton, Humphrey Shotton, and Hywel T. P. Williams. Using social media to detect and locate wildfires. In *EcoMo@ICWSM*, 2016.

[2] Konstantinos-George Thanos, Andrianna Polydouri, Antonios Danelakis, Dimitris Kyriazanos, and Stelios C.A. Thomopoulos. Combined deep learning and traditional nlp approaches for fire burst detection based on twitter posts. In Evon Abu-Taieh, Abdelkrim El Mouatasim, and Issam H. Al Hadid, editors, *Cyberspace*, chapter 6. IntechOpen, Rijeka, 2019. doi: 10.5772/intechopen.85075. URL https://doi.org/10.5772/intechopen.85075.

---

[5] https://nominatim.org/release-docs/latest/api/Overview (accessed 2022-09-25)

[6] https://fireloc.org (accessed 2022-09-28)