

# Dictionary-based Classification of Tweets About Environment

Michela Cameletti<sup>a</sup>, Silvia Fabris<sup>a</sup>, Stephan Schlosser<sup>b</sup>, Daniele Toninelli<sup>a\*</sup>

*a. University of Bergamo; via dei Caniana, 2, 24127, Bergamo, Italy; emails: michela.cameletti@unibg.it; silvia.fabris@unibg.it; daniele.toninelli@unibg.it*

*b. University of Göttingen, Goßlerstraße 19, 37073 Göttingen, Germany; email: stephan.schlosser@sowi.uni-goettingen.de*

*\*Corresponding author: Daniele Toninelli (email: daniele.toninelli@unibg.it).*

## Abstract

In the era of social media, the huge availability of big data such as digital data (e.g. posts sent through social networks or unstructured data scraped from websites) allows to develop new types of research in a wide range of fields. These types of big data are available for low costs and in almost real-time. Nevertheless, their collection and analysis are challenging. This paper proposes an unsupervised dictionary-based method to filter tweets related to a specific topic, i.e. environment. We start from the tweets sent by a selection of Official Social Accounts clearly linked with the subject of interest. Then, we identify a list of expressions (bigrams, trigrams and hashtags) used to set the topic-oriented dictionary. Our approach has some relevant advantages: it attempts to reduce as much as possible the interventions and decisions of the researcher as well as the processing time; it is based mostly on combination of words (instead of single words) in order to ease the identification of tweets concerning the topic of interest; it is not based on a pre-defined dictionary, but it can rather be personalized and generalized to other topics. We test the performance of our method by applying the built dictionary to a sample of more than 3.5 million geolocated tweets posted in Great Britain between January and May 2019. All the criteria used to evaluate the performance highlighted very good performances. In particular, the level of accuracy, of sensitivity and of the F1 score were equal or higher than 98.4%; moreover, also for specificity and precision we obtain excellent levels of performance (around 97,5%), higher than the currently most common methods of selection.

*Keywords:* tweet filtering, big data analysis, dictionary-based selection, dictionary-based search, unsupervised algorithm, text analysis.

*Data and R code:* The datasets generated and analyzed for this paper and the R code used are available in the github repository; link: [https://github.com/silviafabris/Twitter\\_dictionary\\_based\\_classification](https://github.com/silviafabris/Twitter_dictionary_based_classification).

*Authors' contributions:* This is a joint work made by the contribution of all authors. All authors read and approved the final manuscript.

## 1 Introduction

Despite being born as communication tools, nowadays social media and microblogging platforms have become a common source of big data that can be used to develop and/or to support scientific research in a wide range of fields. The collection and the analysis of such data is still an evolving and very promising research field. There are clearly some advantages: for example, data obtained from the Internet are available at lower costs, in shorter times and are easier to be collected than data obtained by means of traditional

surveys. Nevertheless, researchers cannot control the data collection phase: if questionnaire-based surveys ask specific answers connected to particular research purposes, microblogging texts or web scraped data rely on “listening”, collecting and “measuring” what it is available. This means that collected data can be, at least in part, out of focus or too much general (for the chosen research theme) or they could even result misleading. Moreover, since no sampling strategy drives the data collection, the information obtained from microblogging platforms or from the web could be non-representative of the population of interest.

Our research is focusing on data collected from Twitter. Twitter is one of the most relevant and used microblogging platform worldwide, with 321 million active users in February 2019<sup>1</sup>. Using this platform, it is possible to send short messages (up to 280 characters), to retweet them or to like or comment posts shared by other users. Thanks to the wide and still increasing spread of Twitter, a huge (rather than big) amount of data is produced daily. This raw material can be used for developing various types of research on a wide range of fields, such as health (Alessa and Faezipour, 2019; Gesualdo et al., 2013) epidemiology (Ahmed et al. 2019), politics (Budiharto and Meiliana 2018), customer satisfaction (Liu et al. 2017; Hawkins et al. 2016) or well-being and happiness (Mitchell et al. 2013; Baylis et al. 2018). However, retrieving, filtering and processing Twitter data can be extremely challenging. Shared posts can be about personal opinions, ideas, goals and events, but they also include advertisements and news. Moreover, the identification and the selection of tweets regarding a specific topic is also a difficult task. This is mainly due to the completely unstructured nature of this big data as well as to the limited length of posts. There are no precise rules about what to post and how to share information: e.g., users can write plain text or use hashtags to refer to specific topics, to express opinions about an event or to highlight a theme or a fact. Furthermore, news media, government organizations, ONGs, no profit associations, industries and so forth share posts containing information related to their activity, news or advertising.

The purpose of this paper is to propose and test a dictionary-based method to filter tweets concerning a specific topic. This is a first unavoidable step aimed at “making order” among unstructured and undifferentiated big data. Dictionaries are frequently used as base for developing twitter data analyses (Nielsen, 2011). We focus particularly on environment, given its societal importance and its relevance in driving future governmental and cross-national policies. Moreover, as shown in Toninelli and Cameletti (2018) environment is one of the most emerging and important latent dimensions linked to the citizens’ personal well-being. This topic is also particularly challenging, from our point of view, because in all its facets (e.g. climate change, recycling, renewable energy, global warming), it is associated to a language that changes frequently and evolves quickly<sup>2</sup>. For this reason, it is necessary to constantly update the dictionary used to retrieve or select tweets regarding this theme. Consequently, an as-much-as-possible unsupervised and flexible algorithm is needed in order to set up or update a dictionary, whenever necessary.

The approach we propose builds a dictionary starting from a list of keywords obtained by analyzing tweets published by a selected list of Official Social Accounts (OSA) whose activity is strictly related to environment. Once set, we test the performance of this dictionary using a random sample of 3.5 million tweets selected among 54,135,006 tweets posted in Great Britain (GB) between 2019/01/14 and 2019/05/13. These tweets are fully geolocated, given the adopted collection method, based on the “theory of circles”

---

<sup>1</sup> Source: <https://en.wikipedia.org/wiki/Twitter> (latest access: September 3<sup>rd</sup>, 2019).

<sup>2</sup> Source: <https://www.theguardian.com/environment/2019/may/17/why-the-guardian-is-changing-the-language-it-uses-about-the-environment> (latest access: September 3<sup>rd</sup>, 2019).

(Schlosser et al. 2019). As a consequence, it would be possible to study, for example, the spatial variability of the sentiment or the inclination for a certain topic (environment, in our case) across the country's sub-areas.

### 1.1 Added Value

Our algorithm shows competitive features, from the technical point of view, in comparison to what already available in the literature.

First, we propose a flexible method that can be easily applied to any topic of interest. Moreover, the method is also flexible in allowing the researcher to expand as necessary the list of expressions and hashtags included in the dictionary as well as to update and renew the dictionary any time it is needed. This should help in facing the issue of potential frequent changes that can characterize the variable language of topics such as the environment.

Second, our method limits as much as possible the researcher interventions in the process, requesting it mainly in the selection of OSA and in setting few thresholds (e.g. the frequencies used to identify the most common bigrams/trigrams and hashtags).

Third, the algorithm we propose does not need a starting dictionary or a list of predefined keywords, that could be out-of-date or could not perfectly fit with the chosen topic. Thus, in our case the user creates an ad-hoc and personalized dictionary for filtering and selecting tweets linked to a specific theme.

Fourth, the algorithm can be implemented relatively quickly, making it possible to update the dictionary whenever necessary (for example depending upon how much the language of the studied topic is expected to develop and change over time).

Finally, our approach encourages and allows for an "indiscriminate" and massive data collection that (despite being computationally challenging) can be used for any type of purpose and for studying any kind of topic. This, on the one hand, can give rise to computational issues, linked to the necessity to store and manage a potentially huge amount of big data. Nevertheless, on the other hand this represents a highly competitive advantage, mostly if the data collection takes place with reference to specific and reduced time lags: one does not need to run again the full data collection if something was incorrectly set or got wrong during the retrieval phase.

### 1.2 Research Hypothesis

Our method seems to be characterized by some relevant advantages. But does it really and properly work? Is it able to show improvements in the quality level of its outputs, in comparison to other existing methods? Are we really able to filter the amount of available big data selecting what we are really interested in?

The main purpose of this paper is to evaluate the performances of our method in filtering tweets about the specific chosen topic (i.e., the environment).

As a consequence, our **research hypothesis** is: *in addition to the competitive advantages it has, is our method able to enhance the performances of other Tweet filtering methods available in the literature, with respect to the performance indexes we take into account?*

In order to evaluate this, we consider the most commonly used indicators: accuracy, sensitivity, specificity, precision, F1 score (for further details, see sect. 5). Generally, we expect to find support for our hypothesis in our finding. **Accuracy** can be considered an overall measure of the algorithm performances: it focuses on tweets correctly classified as both linked and not linked with the chosen topic. Since our method is based on a specifically-built, up-to-date and personalized dictionary, we expect, generally, to obtain better

performances than the ones obtained by other methods. Nevertheless, in evaluating such an algorithm we could also be mostly interested in measuring the percentage of the selected posts that are really linked to the studied topic. We could measure this by the **Sensitivity**. Higher level of sensitivity reduces the probability of potential biases, in processing our data, or the risk of long further cleaning operations. Our method is specifically built for our purposes. Moreover, it relies on a dictionary that is not based on single words, but on combinations of words (i.e. on bigrams and trigrams), enriched by the addition of the most common hashtags. This could help in obtaining more refined results, in comparison to other algorithms. This feature should also help in reducing the inclusion of non-pertinent tweets and in leading to data of high quality, excluding tweets not linked to the topic chosen (i.e. maximizing the level of **Specificity**) and selecting exclusively tweets really linked to it, among the ones classified as fitting (enhancing the **Precision**). As a consequence of what said before, we expect also to improve the performances in terms of **F1 Score**, obtaining a combined reduction of false positives and false negatives.

All the studied indexes confirm the very good performance obtained by our algorithm. This represents an important contribution in enhancing the literature about how to select/filter messages sent through social media about a specific topic.

This paper is structured as follows: in Sect. 2 we frame our paper within the existing literature and we review the main works related to the classification and the tweets filtering; Sect. 3 describes the analyzed data; Sect. 4 introduces the algorithm set to build the dictionary and to filter tweets; Sect. 5 discusses the results of the quality metrics we used for testing our algorithm using the GB tweets; Sect. 6 discuss our main findings and provides ideas for further studies.

## 2 Literature Review

In the framework of big data research and analytics, works on data retrieved from social networks is taking a relevant role. In particular, among the Twitter-connected literature, a big part is reserved for topic detection studies (see e.g. Ashgari et al. 2020, Snyder et al. 2020, Mottaghinia et al. 2020). These are mainly aimed at classifying tweets into different topics, known or not known a priori. In this paper we consider the case of a pre-defined topic, i.e. environment. Starting from known categories, the literature have been proposed two main methodological approaches: dictionary-based and supervised methods (Grimmer and Stewart 2013).

The first approach filters tweets about a specific topic by using a set of keywords defining a dictionary. For example, Cody et al. (2015) explore climate change sentiment by selecting tweets containing at first the word “climate”, and related expressions such as “global warming”, “climate realist”, “climate change” and “anthropogenic global warming”. Similarly, in order to collect tweets related to ecological crisis, Kozlowski et al. (2020) use both generic (e.g. flood, storm) and specific keywords (e.g., crisis names and types, such as Irma hurricane). It is also possible to perform the selection of tweets by considering specific hashtags related to the topic of interest. In this regards, Reyes-Menendez et al. (2018) study the opinion about environment by selecting posts containing the hashtag #WorldEnvironmentDay. A similar approach has been used by Pruss et al. (2019). For filtering tweets regarding the Zika infection, they first use the keywords “zika” and “ZIKAV” to find posts; then, they apply a topic model like the Latent Dirichlet Allocation to find the most popular sub-topics (e.g. environmental concerns, vaccination, emergency declaration). The use of the hashtags to select tweets is also adopted in Harb et al. (2020) to identify tweets related to mass violent events.

These strategies based on keywords and hashtags are very simple and easy to be applied; nevertheless, they do not guarantee to identify all the tweets connected to the specific argument. At the same time a large set of keywords or hashtags can lead to the inclusion of posts not strictly related to the topic of interest.

Supervised methods, instead, require the human intervention to create a (large) dataset of labeled tweets (i.e. classified into a predefined set of categories) that will be used for training (or supervising) a statistical model or a machine learning algorithm. This estimated classifier is then adopted for predicting the category of a new tweet. This approach is used for example by Frenda et al. (2019), who adopt a common method in machine learning as Support Vector Machine (SVM) to automatically detect sexist and misogyny on Twitter. In particular, this study uses, as training dataset, the freely available corpora known as “Automatic Misogyny Identification – IBEREVAL 2018”<sup>3</sup> and considers, as input variables, some lexical and stylistic features of the post. The accuracy of such a classifier is equal to 76%. Similarly, Foucault and Courtin (2016) combine SVM and a Naïve Bayes method to classify tweets sent from French institutions into four communication categories (sharing experience, promoting participation, interacting with the community, and promoting-informing about the institution). Each tweet is represented by 18 features derived from metadata information, punctuation marks, tweet-specific characteristics (e.g. use of hashtags and emoticons) and lexical features. By means of cross-validation, the classification performance is evaluated obtaining a value of the F1 index<sup>4</sup> equal to 72%. This kind of approach is also known as “feature-based modeling”, because it requires to extract textual features from the tweets to be successively provided as input to the machine learning algorithms. A different approach is applied in Ghafarian and Yazdi (2020) for the identification of informative posts during crisis episodes (e.g. earthquakes, floods, etc.). In particular, they consider each tweet as a distribution from which a sample of words is drawn and adopt Support Measure Machine, which is an extension of SVM for distributional data. They apply their method to 19 different crisis datasets and obtain values of the F1 index between 70.7% and 87.5%.

An extension of this supervised strategy is represented by deep learning methods. They employ neural networks and usually outperform feature-based models. Nizzoli et al. (2019), for example, analyze extremist propaganda and try to identify pro-ISIS tweets. In particular, they show that a Recurrent Neural Network (RNN) with pre-trained word embedding is able to reach a F1 score of 0.9. Stowe et al. (2018) adopt SVM and linguistic features to identify tweets related to hurricane events. As a comparison, they implement also two deep learning methods (Multi-layered Perceptron and Convolutional Neural Network); they find that the Multi-layered Perceptron performs better with a F1 score equal to 0.83. They also run feature-based algorithms (logistic regression, SVM and Naïve Bayes algorithm) with linguistic, temporal and geospatial features to predict people behavior during hurricane events. In this case SVM provides the best performance with F1 values included between 0.47 and 0.79 according to the considered features. Also Grzeża et al. (2020) adopt deep learning methods for the automatic classification of tweets related to alcohol consumption. In particular, they propose an ensemble of two classifiers based on distributional semantics and Convolution Neural Network (CNN), respectively. By using five different datasets, they obtain values of the F1 index between 79.7% and 94.6%. Pamungkas et al. (2020) use SVM (with word n-grams as features) and RNN to detect misogyny in tweets. By considering datasets written in three different languages (English, Italian, and Spanish), they obtain values of the F1 score ranging from 39.2% to 89.1%. The problem of classifying

---

<sup>3</sup> Source: <https://amiibereval2018.wordpress.com/> (latest access: September 3<sup>rd</sup>, 2019).

<sup>4</sup> The F1 score is a performance index depending on precision and recall (see Sect. 5 for its definition).

health-related tweets for the early detection of disease outbreaks is developed in Şerban (2019) by using CNN and RNN. Their best classifier gives an F1 of 85.2%.

The method we propose in this paper is unsupervised and dictionary-based. Differently from Cody et al. (2015) and Pruss et al. (2019), it builds a dictionary that includes the most common bigrams, trigrams and hashtags about environment. A bigram (trigram) is a pair (triple) of consecutive words commonly associated one to each other. Starting from bigrams, trigrams and hashtags about a specific topic, our method does not need a starting predefined set of keywords. Only a list of OSA has to be provided in advance and a limited number of human checks are needed in order to avoid the inclusion of too general keywords that would lead to the selection of tweets not strictly pertaining to environment (or to the chosen topic). Thus, our approach minimizes the amount of required human work, because it doesn't need the set of labeled tweets for training (as required by supervised approaches) or a predefined set of keywords (that could be too much general or not completely focused on the studied topic). At the same time, thanks to the arbitrary selection of the OSA and to the possibility of reviewing step by step the dictionary creation, it is very flexible and could be applied to and personalized for any topic of interest.

### 3 Data Collection and Preliminary Cleaning

In the following subsections we describe the two datasets used for the analysis. The first one (Sect. 3.1) is composed by a sample of tweets posted by OSA related to the analyzed topic, environment. Starting from these data, the algorithm sets up the dictionary. The latter is then applied to the second dataset (Sect. 3.2), composed by the tweets posted in GB between 2019/01/14 and 2019/05/13. The algorithm has been implemented using the R software<sup>5</sup>.

#### 3.1 Tweets from OSA

The general idea behind the tweet selection is that posts speaking about the same topic should be similar and different from tweets related to other themes. As a consequence, tweets pertaining to a certain topic should generally include similar words or combination of words. Our work aims at detecting and studying posts about environment. For this purpose, our preliminary objective is to set up a dictionary including the most common and relevant keywords related to such a topic. As first step, we identified 12 OSA linked to environment. In particular, we chose verified accounts<sup>6</sup> (or profiles that have at least 10,000 followers) belonging to no-profit associations, research institutes and intergovernmental organizations whose activity is related to environment<sup>7</sup>. The OSA selection is an arbitrary phase of the algorithm. The chosen accounts are selected because of their popularity and with the aim of covering all the possible aspects of environment (e.g. climate change, plastic pollution, nature protection). Note that, as it will be described in Sect. 4.2, the OSA choice can cause effects on the final dictionary.

---

<sup>5</sup> The code and the data are available at the following link:

[https://github.com/silviafabris/Twitter\\_dictionary\\_based\\_classification](https://github.com/silviafabris/Twitter_dictionary_based_classification)

<sup>6</sup> <https://help.twitter.com/en/managing-your-account/about-twitter-verified-accounts> (latest access: September 3<sup>rd</sup>, 2019).

<sup>7</sup> @climateprogress (Climate Progress), @ClimateReality (Climate Reality), @friends\_earth (Friends of the Earth), @Greenpeace (Greenpeace), @GreenpeaceUK (Greenpeace UK), @LessPlasticUK (Less Plastic), @PlasticPollutes (Plastic Pollutes), @UNEnvironment (UN Environment Programme), @UNFCCC (UN Climate Change), @World\_Wildlife (World Wildlife Fund), @WWF (WWF), @WWFScotland (WWF Scotland).

For each account, we retrieved all the most recent posted or retweeted tweets that Twitter leads us to download up to 2019/05/10. Among the obtained 38,611 tweets, we kept exclusively posts written in English. Then, we cleaned their corpus by removing url links, html codes, non-ascii and special characters, but we kept hashtags. This list of cleaned tweets is our first dataset. In Sect. 4.1 we analyze this dataset, in order to detect the most recurrent expressions (i.e. bigrams, trigrams) and hashtags used by the considered OSA.

### 3.2 Tweets from GB

A second dataset is used in order to test and apply the dictionary built with the selected keywords starting from the first dataset introduced in sect. 3.1. This second database includes all the tweets sent in GB from 2019/01/14 to 2019/05/13 (i.e. for a total of 120 days). The tweets are collected through the “theory of circles” method, described in Schlosser et al. (2019) and further tested in Schlosser et al. (2020). Generally, just 1-2% of tweets contains GPS coordinates<sup>8</sup>. Nonetheless, the “circle approach” allows us to geolocate all tweets directly in the collection phase, associating each post to one of the NUTS-1<sup>9</sup> sub-areas covering GB (see Figure 3, left).

After having preliminary removed messages sent by bots<sup>10</sup>, we obtained 54,135,006 tweets, that corresponds to an average of 4,921,364 tweets for each NUTS area.

The next step consisted in cleaning the tweets’ corpus. In doing so, we tried to keep as much information as possible by replacing htmls, emojis and slangs with equivalent-meaning expressions. For example, the Unicode character “\U0001f602”, corresponding to the emoji 😂, was translated with “face with tears of joy”. Note that we kept hashtags in the text because they are crucial in detecting tweets related to the specific topic. The cleaned data of this second dataset are then used in order to compute some performance indexes for our dictionary (see Sect. 5.1).

## 4 Methods

After having cleaned the posts sent by the selected OSA (Sect. 3.1), we use them in order to set up the dictionary (the first dataset, introduced in 3.1). The steps of the algorithm for the dictionary definition are explained in Sect. 4.1 and summarized in the flow charts of Figure 1 and 2. The final dictionary is composed by a set of selected bigrams, trigrams and hashtags and is applied to the full set of GB tweets described in Sect. 3.2 (the second dataset) in order to select tweets regarding the chosen topic.

### 4.1 Selection and cleaning of bigrams, trigrams and hashtags from OSA tweets

Given the tweets collected from the selected OSA and preprocessed (see Sect. 3.1), we produce the list of all bigrams and trigrams (also named *expressions* in the following) with the corresponding frequencies (see Table 1). This represents the starting point of the dictionary creation (step *a* in Figure 1). Expressions which do not appear frequently are usually not related to the topic or are too general to be included in the final dictionary. For this reason, in order to select the most pertaining bigrams and trigrams (that will be later used to define the dictionary), some additional steps are required. After having looked at the full list of bigrams

<sup>8</sup> <https://developer.twitter.com/en/docs/tutorials/tweet-geo-metadata.html> (last access on August 2nd, 2019)

<sup>9</sup> Source: <https://ec.europa.eu/eurostat/web/nuts/background> (latest access: September 3<sup>rd</sup>, 2019).

<sup>10</sup> A bot is an automated program which interacts automatically on the social network.



and trigrams not containing stop words<sup>11</sup>, it is possible to proceed directly with the definition of some thresholds for the frequencies in order to exclude expressions which do not occur often (step *c* in Figure 1). In our case, we exclude bigrams which appear less than 65 times and trigrams posted less than 35 times (see the non-grey expressions in Table 1). However, the resulting list of expressions could still include bigrams and trigrams not related to environment (e.g. common or general expressions, country and state names). Consequently, an additional cleaning stage is required (step *b* in Figure 1). In particular, the algorithm proposes the following possibilities:

- i. a *standard* cleaning (step *d* in Figure 1) to be performed after the choice of the frequency thresholds (step *c* in Figure 1): the user reviews the expressions selected after applying the thresholds and remove any expressions not strictly related to the topic;
- ii. an optional *extra* cleaning stage before the choice of the thresholds (step *b* in Figure 1). The aim is to remove, from the list of selected OSA expressions, some common terms which are widely used in Twitter and very likely not related to environment. Even if this step is optional, we highly suggest to use it, because it reduces the standard review process performed at step *d* in Figure 1. The extra cleaning considers the full set of GB tweets described in Sect. 3.2 to identify the list of *general expressions*, i.e. popular bigrams and trigrams (step *a* in Figure 2). These recurrent expressions are used to remove from the OSA bigrams and trigrams list general expressions such as “trump administration”, “taking action”, “million people”. It is important to note that this procedure can be performed by using the full set of GB tweets or a smaller sample, in order to reduce the computational time. Our empirical experience with our case study demonstrates that the final dictionary does not change considerably by using different samples or the complete dataset of GB tweets. For this reason, we decided to use a random sample of 3.5 million tweets collected between March 10<sup>th</sup> and May 13<sup>th</sup>, 2019. We arbitrary decided for a very low threshold: in the list of general expressions we take into account just bigrams and trigrams tweeted at least 20 times. This way, we obtained 30,656 general expressions (step *a* in Figure 2). However, this vector of recurrent bigrams and trigrams may contain expressions linked to environment, such as “climate change”, which we do not obviously want to be part of the list (otherwise they will be not included in the dictionary). Thus, a review of the list of general expressions is necessary. This can be done by adopting one of the following two approaches:
  - a. *user-based approach* (step *b* in Figure 2): the user examines all the general expressions one by one and remove the ones related to environment;
  - b. *list-based approach* (step *c* in Figure 2): in this case we assume that a set of expressions related to environment is available (prepared ad hoc by the researcher or taken from existing dictionaries). The two lists will be matched and the environment-related expressions will be removed from the set of general expressions.

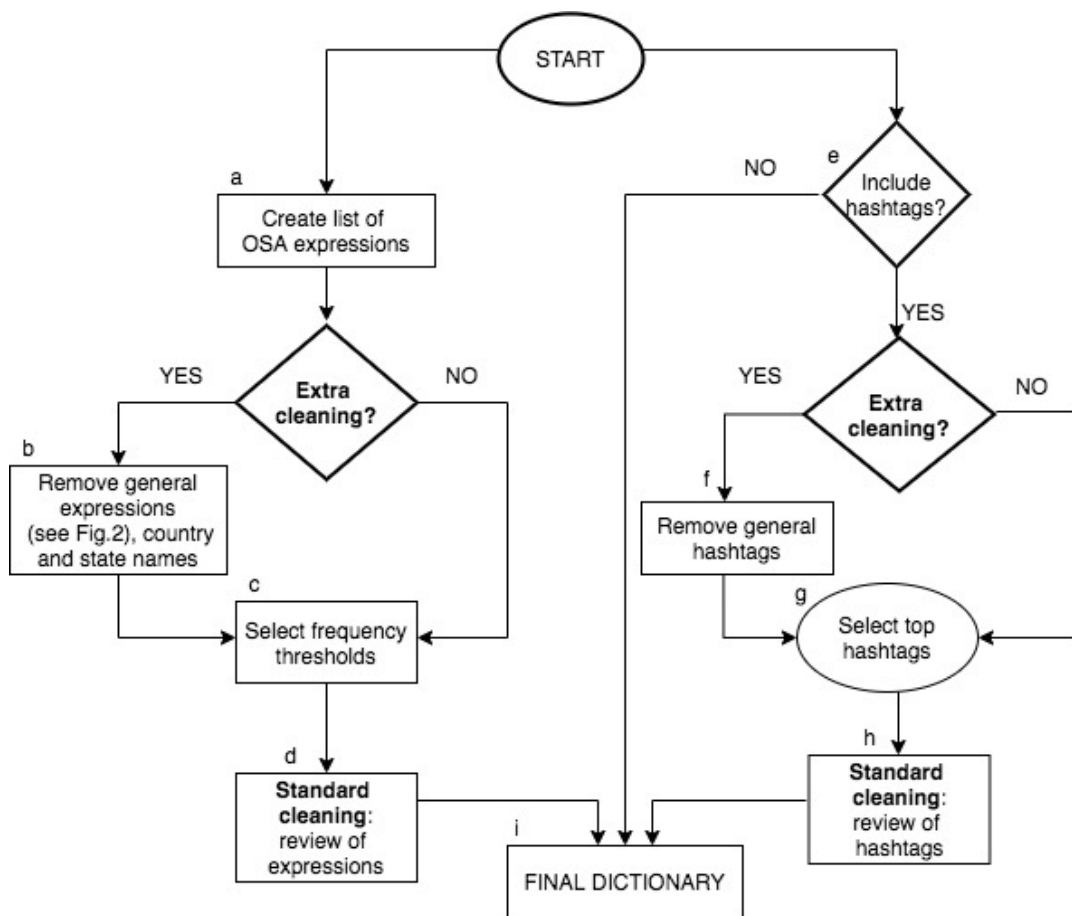
---

<sup>11</sup> Stop words are the most common words used in a language (such as for example “the”, “a”, “an”, “in”). In this case the list of stop words is given by three different lexicons (“onix”, “SMART” and “snowball”).



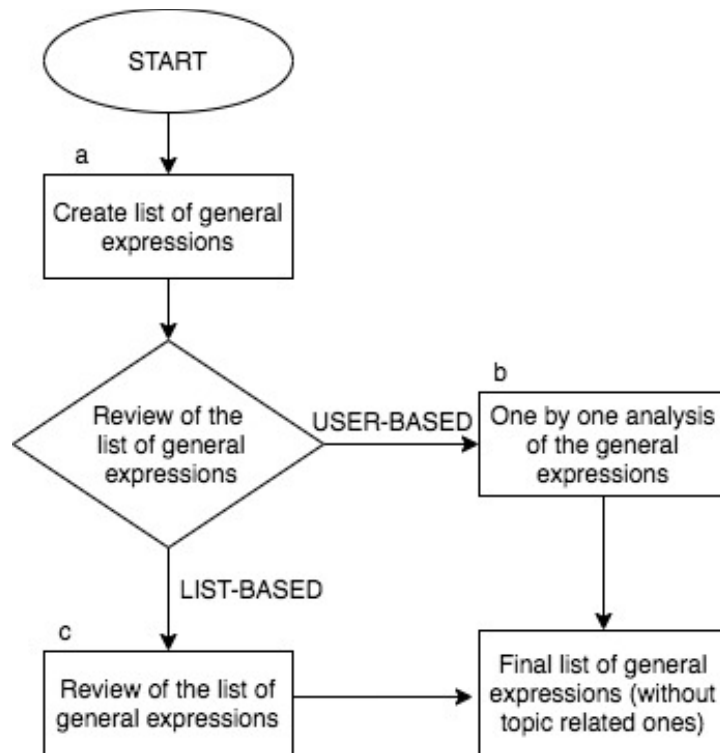
## Dictionary-based Classification of Tweets About Environment

In this research we adopted the list-based approach by considering a list of 49 expressions<sup>12</sup> prepared specifically for our case study. After removing from the list of general expressions the terms related to environment, the resulting vector is composed by 30,632 expressions. It is then possible to proceed with the extra cleaning of the OSA bigrams and trigrams by removing all the terms included in the set of 30,632 common expressions (step *b* in Figure 1). Moreover, in this same step, all the bigrams and trigrams that contain country names and USA state names are removed. The extra cleaning step removed 7 bigrams and no trigrams (see the red expressions in Table 1). Finally, after the *extra* cleaning, the *standard* cleaning (step *d* in Figure 1) is used to review the new list of OSA expressions in order to exclude other terms not related to the studied topic, such as “start donating” or “coral reefs” (see blue expressions in Table 1). For our application this standard review step removed 10 bigrams and 1 trigram. As result, we obtain the final list of bigrams and trigrams related to the topic.



**Figure 1** Algorithm for creating the final dictionary

<sup>12</sup> air clean, air pollution, air quality, carbon emissions, carbon pollution, clean air, clean energy, climate action, climate change, climate conference, climate crisis, climate reality, climate science, climate solutions, coal ash, coal plants, coal-fired power, conference cop, environmental laws, extreme weather, food waste, fossil fuel, fossil fuels, fuel industry, gas drilling, gas emissions, gas industry, global climate, global temperatures, global warming, greenhouse gas, healthy environment, offshore drilling, palm oil, paris agreement, plastic bags, plastic bottles, plastic packaging, plastic pollution, plastic straws, plastic waste, renewable energy, single-use plastic, single-use plastics, tar sands, toxic chemicals, toxic pesticides, warming world, weather events.



**Figure 2** Extra cleaning phase (step *b* of Figure 1)

After the previous phase, the algorithm gives the possibility to enrich the dictionary with hashtags (step *e* of Figure 1), which are normally used by Twitter users to identify and categorize tweets. In our case we decided to include hashtags because they can be extremely useful in filtering a tweet linked to environment. First of all, we analyze the hashtags used by the selected OSA accounts (see Table 2). They represent either popular themes on Twitter (i.e. trends), slogans used in the OSA description (e.g. #UseLessPlastic for the @LessPlasticUK account) or they are created by OSA for particular international events (e.g. #PlasticFreeFriday). As done previously for OSA expressions, we need to remove hashtags too general (such as “#nature”) and referring to countries or to states. Even in this case, it is possible to adopt the *standard* review only (step *h* of Figure 1), or to also use an *extra* cleaning (step *f* of Figure 1). For this purpose, we create a list of general (i.e. popular) hashtags by using a sample of the GB tweets (as described in Sect. 3.2). After removing the general hashtags (see red terms in Table 2), we selected the 60 most popular OSA hashtags (step *g* of Figure 1). Finally, a last standard review (step *h* of Figure 1) is implemented with the aim of excluding the hashtags that are too generic to be part of the dictionary (such as “#climate” or acronyms like “#dyk”, i.e. “Do You Know”). These excluded hashtags are reported in blue in Table 2.

The final dictionary is composed by 35 OSA expressions (listed in Table 1) and by 52 hashtags (see Table 2). We apply this dictionary to a sample of 3.5 million tweets randomly selected over 54 million tweets collected in GB in order to select only tweets that contain at least one expression included in the final dictionary. As a result, we obtained 107,176 tweets related to environment.

## 4.2 The effect of the number of considered OSA

The number of OSA considered for defining the dictionary can be arbitrary chosen by the user, as described in Sect. 4.1. In order to assess whether there is an effect of the number of selected OSA on the final list of expressions defining the dictionary, we compare the list of bigrams/trigrams obtained using the 12 OSA described in Sect. 4.1 (named LIST 1, in the following) with the one obtained by using 22 OSA (the previous 12 OSA plus 10 new ones<sup>13</sup>, named LIST 2). With LIST 2 the number of OSA tweets increases from 38,611 to 57,029. Moreover, by applying the extra cleaning and keeping the same thresholds set in Sect. 4.1, we obtain a larger final list of expressions composed by 74 bigrams/trigrams (instead of 35). This means that increasing the numbers of OSA leads to a larger set of expressions but also to a more demanding standard cleaning step. As a consequence, we suggest to keep the number of OSA between 10 and 15 in order to avoid unnecessary cleaning to remove expressions which are too generic and not strictly related to the topic of interest and are highly time consuming.

Moreover, by comparing the final expression list obtained with LIST 1 and LIST 2, it can be observed that: *i)* there is just 1 expression in LIST 1 which does not appear in LIST 2; *ii)* there are 40 expressions contained in LIST 2 which are not included in LIST 1; *iii)* considering the 20 most recurrent expressions, the two lists differ by just 6 terms (3 are contained in LIST1 and not in LIST2 and, contrarily, 3 are included in LIST2 and not in LIST 1). On the basis of these results, we can conclude that, even if LIST2 gives rise to a larger set of expressions, by looking at the most frequent terms the two lists are almost equal.

**Table 1** OSA bigrams and trigrams and corresponding frequencies.

Bigram / Trigram	Freq.	Bigram / Trigram	Freq.
climate change	1357	leadership corps	81
donating tweet	547	million people	81
start donating	547	reality leadership	80
tweet unsubscribe	547	climate conference	79
plastic pollution	460	singleuse plastics	79
climate crisis	405	climatechange conference	77
climate action	302	extreme weather	77
air pollution	283	air quality	74
palm oil	253	offshore drilling	70
climate reality	246	coral reefs	68
singleuse plastic	245	paris agreement	68
renewable energy	213	tar sands	68
plastic waste	202	food waste	66
fossil fuel	183	marine life	66

<sup>13</sup> @climateprogress, @ClimateReality, @friends\_earth, @Greenpeace, @GreenpeaceUK, @LessPlasticUK, @PlasticPollutes, @UNEnvironment, @UNFCCC, @World\_Wildlife, @WWF, @WWFScotland, @NRDC, @nature\_org, @EnvDefenseFund, @Earthjustice, @foe\_us, @guardianeco, @HuffPostGreen, @insideclimate, @PlanetGreen, @ClimateCentral.

Dictionary-based Classification of Tweets About Environment

clean energy	180	world leaders	66
fossil fuels	176	antarctic ocean	65
greenhouse gas	147	natural gas	65
plastic packaging	128	deposit return	63
uselessplastic lesssoceanplastic	125	ocean plastic	63
plastic straws	124	uk government	63
connect earth	119	david attenborough	62
trump administration	118	barrier reef	61
gas emissions	117	sea ice	61
global warming	107	raise awareness	60
carbon emissions	99	[more excluded bigrams]	
conference cop	97	antarctic ocean sanctuary	43
global climate	94	support balloon releases	43
taking action	90	rising global temperatures	29
plastic bags	89	arctic sea ice	28
scott pruit	89	conference sb bonn	26
send thinkprogress	86	drasticonplastic timer challenge	24
human health	84	action summit gas	22
national park	84	exposing white nationalism	22
clean air	82	fashioned po box	22
plastic bottles	82	mobile calendar wallpaper	22

*Note:* Gray: excluded (frequencies lower than 65/35 for bigrams/trigrams). Red: removed by the extra cleaning (step a, Figure 1); blue: removed by the manual cleaning (step d, Figure 1).

**Table 2** OSA hashtags and corresponding frequencies.

Hashtag	Freq.	Hashtag	Freq.
#climatechange	1530	#youngchamps	86
#climateaction	861	#renewables	85
#cop	759	#earthday	83
#plasticpollutes	629	#refusesingleuse	82
#endangeredemoji	566	#renewableenergy	80
#parisagreement	461	#reuse	79
#plasticpollution	389	#oceanplastic	78
#earthhour	375	#passonplastic	77
#uselessplastic	337	#biodiversity	75
#climate	301	#beatairpollution	74
#beatplasticpollution	277	#climateambition	72

## Dictionary-based Classification of Tweets About Environment

#earthhourscotland	244	#nature	70
#breakfreefromplastic	237	#bethechange	69
#plastic	207	#talanoa	65
#fracking	193	#nothirdrunway	63
#cleanseas	191	#fridaysforfuture	62
#beatpollution	178	#zerowaste	61
#globalgoals	174	#africaclimateweek	61
#pandahugs	161	#cleanenergy	60
#lessoceanplastic	154	#beachclean	60
#solvedifferent	150	#wildforlife	60
#plasticfree	149	#breathelife	59
#protectantarctic	136	#fightforyourworld	58
#lessplastic	132	#singleuseplastic	57
#connect	129	#climatebreakdown	56
#airpollution	126	#climatechangebill	56
#climateemergency	125	#renewable	56
#worldenvironmentday	125	#solar	55
#dyk	120	#climatestrike	55
#endoceanplastics	111	#atlanta	53
#gcas	97	#oneplanet	52
#actonclimate	95	#blueplanet	51
#promisefortheplanet	94	#climatehope	50
#sb	89	#worldwildlifeday	50
#dropdirtypalmoil	88	#bees	50
#drasticonplastic	87	#reusable	50

*Note:* Gray: excluded (not belonging to the top 60). Red: removed by the extra common cleaning (step f, Figure 1); blue: removed by the standard review (step h, Figure 1).

## 5 Results: Dictionary Performance

In order to evaluate the performance of the dictionary-based filtering, we randomly choose 600 tweets selected and 600 not selected by the algorithm (i.e. classified as not linked to “environment”). Then, we manually classify these posts into two categories: “related” and “non-related” to environment. This allows us to compute the following relevant quantities, which can be collected in the confusion matrix reported in Table 3:

- a. number of **true positive** (TP), i.e. number of tweets correctly classified by the algorithm as related to environment;
- b. number of **false positive** (FP), i.e. number of tweets wrongly classified by the algorithm as related to environment;

Dictionary-based Classification of Tweets About Environment

- c. number of **true negative** (TN), i.e. number of tweets correctly excluded by the algorithm because not linked to environment;
- d. number of **false negative** (FN), i.e. number of tweets wrongly classified by the algorithm as not pertaining the environment.

**Table 3** Confusion matrix.

		True category	
		Not related	Related
Predicted category	Not related	TN	FN
	Related	FP	TP

The algorithm performance has been evaluated through the following indexes, based on the confusion matrix shown in Table 3:

- a. **accuracy** (AC): proportion of tweets correctly classified by the algorithm on the total number of tweets processed:

$$AC = \frac{TP + TN}{TP + TN + FP + FN};$$

- b. **sensitivity** (SE): proportion of tweets pertaining to the argument that are correctly filtered by the algorithm:

$$SE = \frac{TP}{TP + FN};$$

- c. **specificity** (SP): proportion of tweets not related to environment which are correctly excluded by the algorithm:

$$SP = \frac{TN}{TN + FP};$$

- d. **precision** (PR), i.e. the proportion of tweets truly related to the topic among tweets classified by the algorithm as pertaining the chosen topic:

$$PR = \frac{TP}{TP + FP};$$

- e. **F<sub>1</sub> score** (F1) defined as a function of PR and SE and given by:

$$F_1 = 2 \frac{PR \cdot SE}{PR + SE}.$$

For all the indicators the range is between 0 and 1 and the “the higher the better” rule holds. Results, here, are expressed in percentage terms.

AC can be used as a first overall measure to evaluate the classification algorithm performance, taking into account tweets correctly classified on the total number of posts. In particular, following this criterion, our method is able to classify correctly 98.42% of the total number of tweets. This measure of performance is quite high, if compared, for example, with the accuracy obtained by the Automatic Misogyny Identification – IBEREVAL 2018 classifier (equal to 76%). However, AC, used by itself, can be misleading, especially when there is a severe class imbalance in the classification problem.

A more informative evaluation is obtained by using AC together with SE (also known as recall), which represents the ability of the algorithm in correctly selecting all the tweets concerning environment on the total number of tweets linked to the topic. The SE value is even higher than the AC measure: 99.32% of the tweets truly related to environment are identified as relevant by the algorithm. Just 0.68% of the tweets not

linked to environment are wrongly selected by our algorithm, that represents a very small percentage. However, even a high value of SE could hide a problematic situation. This happens when, for example, the number of false positives is high.

In such cases, it is necessary to consider an additional measure as PR. The latter expresses the proportion of tweets truly talking about environment (TP) among the tweets classified by the algorithm as related to the topic (FP+TP). Our algorithm of classification leads to excellent performances also from this point of view: 97.5% of tweets classified as linked to environment are actually speaking about this topic.

Both SE and PR are independent of the number of true negatives. Thus, the evaluation can be improved by taking into account SP in order to evaluate the percentage of tweet not linked to the environment that are correctly excluded by the algorithm: this percentage is equal to 97.5%. A low value of SP would mean that the algorithm has a high rate of false positive. Nevertheless, both considering tweets linked and not-linked to the environment, the algorithm we propose is able to find out a quite high percentage of correctly classified tweets. This means that the capability of our method is very well balanced for both the categories of tweets (linked and not linked with the topic, i.e. TP and TN), at least in our case study.

Given the usual trade-off that exists between SE and PR (i.e. when SE increases, PR decreases and vice versa), it is suggested to combine the two measures in an overall index represented by the F1 score. This is defined as the harmonic mean of SE and PR. As a result, we obtain a value of F<sub>1</sub> score equal to 98.40%. This is very satisfactory and shows that the algorithm has a low rate of false positives and false negatives; this means that we are able to correctly identify relevant messages and, at the same time, do not include in our analysis non-pertinent tweets.

Table 4 summarizes all the performance measures computed for our method.

**Table 4** Performance indexes values (in percentages) for the environment dictionary-based algorithm.

AC	SE	SP	PR	F1
98.42	99.32	97.55	97.50	98.40

All the scores reported in Table 4 are very close to 100 and denote, generally (i.e. from any of the considered point of view), very good performances of the algorithm.

Evaluating all the indexes together, we notice that our dictionary-based algorithm performs very well. For example, it outperforms the supervised feature-based models of Frenda et al. (2019) and of Foucault et al. (2016) which report an AC of 76% and a F<sub>1</sub> score equal to 72%, respectively. In terms of AC, our approach also outperforms the levels obtained by Samuel et al. (2020) applying two machine learning classification methods to classify Coronavirus Tweets. For short tweets, the obtained accuracy was 91% using the Naïve Bayes method and 74% using the logistic regression classification method (whereas the performances obtained for longer tweets are relatively weaker).

Edo-Osagie et al. (2019) applied semi-supervised classification techniques as well as alternative techniques to popular deep learning approaches, including special features such as emojis in order to improve the classification performances. Using the MLP, according to the authors the best fully-supervised approach (in terms of the F1 and F2 scores), they achieved an overall AC equal to 95.5%. The best semi-supervised approach led to an AC equal to 95.5%. The F1 scores obtained by the fully-supervised and by the semi-supervised algorithms were equal to 91.0% and to 91.2%, respectively.



Levels similar to the ones obtained by using our methods were reached by Kudugunta and Ferrara (2018), but for more general purposes than ours, that is for separating bots from human tweets. In this case the authors use a deep neural network based on contextual long short-term memory architecture that exploits both content and metadata to detect bots. From single tweets, the research achieved a high accuracy ( $> 96\%$ ), increased to  $>99\%$  in the case of account-level bot detection.

At the same time our method performs more similarly (despite anyway outperforming such methods of more than 8 percentage points) to the computationally intensive deep learning methods implemented in Stowe et al. (2018) and Nizzoli et al. (2019), which obtain values of the  $F_1$  score equal to 83% and 90%, respectively. Additionally, the computational complexity and the load for the researcher, using our method, is noticeably reduced, in comparison to such intensive deep learning-based approaches.

## 6 Discussion and Conclusions

In this paper we propose an unsupervised dictionary-based algorithm for dealing with huge amount of big data. In particular, our scope is filtering tweets concerning a specific topic: environment. Differently from supervised methods, our approach does not require to have a set of labeled tweets, an operation which is particularly expensive in terms of human work necessary to classify the tweets for creating the training set of data. For this reason, our approach is easier to be implemented. However, environment is a general and wide topic of discussion if compared, for example, with specific events such as earthquakes or other natural disasters. For this reason, the selection of tweets is more difficult because it cannot be performed by using specific keywords strictly referring to the name of the event as done by Kozłowski et al. (2020). When the researcher has to choose the set of keywords to include in the dictionary, as requested by unsupervised algorithms, there is the risk of omitting some important expressions, being the topic so multifaceted. Our approach deals with this criticality by means of expressions adopted by OSA to define the initial list of keywords. As discussed in Sect. 4.1 and depicted in Figure 1, the algorithm follows different potentially iterative steps. In particular, we consider two review phases (standard and extra) both applied to the expressions and to the hashtags as well. The extra cleaning step is useful to reduce the final standard review of expressions, but it is not mandatory and users can proceed with the standard cleaning only. This is a first advantage of our method, in terms of flexibility: the researcher can decide whether (or how frequently) to implement a deeper cleaning phase according to several factor (e.g., how frequently the base language for the topic under study changes, how much time passed after the previous analysis, how much time and resources are available for the processing, and so forth). However, generally we suggest to adopt the standard approach only the first time the user creates a dictionary and to prefer the double review (standard and extra) otherwise, in order to simplify the dictionary creation procedure. In particular, if a list of environment related expressions is already available (i.e. we can implement the list-based approach of Figure 2), our suggestion is to use this big set of general expressions as a base for the extra cleaning. On the contrary, if such a list is not available or is outdated (i.e. the user-based approach of Figure 2), it is preferable to limit the number of general terms that later should be manually checked in order to avoid including environment-related expressions. In any case these two different strategies (list-based vs user-based approach) will not change significantly the final dictionary results, as proved by our results (see sect. 4.2).

Moreover, the flexibility of our methods is even more enhanced by the fact that the thresholds set to select the bigrams and trigrams keywords and the list of hashtags to be included in the dictionary can be increased

or lowered, if needed. These settings can be defined according to the study background, circumstances, resources and can be changed or adapted taking into account, for example, if we are developing a first run rather than further waves of the same analyses. The approach we propose for setting the thresholds is even less prone to the arbitrary of a human decision than it could seem. By visualizing the entire list of bigrams, trigrams and hashtags, it should become clear that at some point the expressions start to be unrelated to the argument of interest. This makes the identification of the correct threshold a relatively easy and a not-so-arbitrary task. Some keywords, such as “coral reefs”, which can appear into the final dictionary, can be relevant to the phenomenon object of study, but at the same time can lead to the inclusion of misleading tweets. We suggest to keep exclusively expressions strictly related to the environment topic, trying, generally, to be more restrictive than inclusive.

With respect to our research question, as shown and discussed in Sect. 5, the indexes obtained for our method highlight very good performances. Both AC and  $F_1$  scores are higher than 98%; but also evaluating the capability of identifying TP and TN the performances are excellent (97.5%). Thus, our method seems to outperform the main algorithms recently proposed in the literature, while being at the same time extremely convenient from the computational point of view (i.e. the running time is extremely lower if compared with deep learning models).

This is not the only advantage of our algorithm. It is extremely flexible, since it can be applied in any type of field dealing with textual big data and studying any type of topic that includes any analysis based on messages sent by using a social media. Moreover, we propose a method that allows the researcher to include a pre-set dictionary, in order to create an own personalized one or allows to integrate both approaches. Being able to update or “increase the size” of a dictionary, when desired, is a valuable advantage, mostly if the topic studied is linked to a language quickly changing and/or to events that drive and characterize the citizen perception of the topic itself. Finally, differently from other dictionary-based methods proposing as starting point a list of single keywords or terms, we propose the use of bigrams and trigrams; this choice reduces the error of misclassification related to the use of single words. At the same time, the inclusion of hashtags and the “translation” of emoticons or web links allow to keep connected the dictionary creation with the most recent updates and events linked to the topic under study.

Our research is currently affected by some limits, that can suggest ideas for further research. The unsupervised method we propose is still not fully automatic, because it requires to set the thresholds for selecting the hashtags and bigram/trigrams and a final manual removal of keywords not strictly related to the topic of interest is needed. Nevertheless, this could also be seen as an advantage, because it makes the algorithm flexible and modifiable according to the users’ requirements or preferences. Moreover, the list of OSA used as base for our algorithm affects the obtained dictionary and this phenomenon could be more relevant in other field or studying other topics. Thus, other less arbitrary criteria about which and how many OSA to select as starting point can be proposed and tested.

Other ideas for further research include deeper and more general studies aimed, for example, to check the degree of generalization of our methods in other contexts (i.e. varying topics and/or countries/languages).

As future research, we intend to increase the number of keywords included in the dictionary by means of a periodic analysis of the OSA accounts. In fact, environment is a really sensitive, discussed and extremely trending topic, which can vary frequently. For these reasons, it is crucial to update the dictionary on a regular basis, in order to capture new trends, impact of events, and the consequent changes of people opinions.

## Dictionary-based Classification of Tweets About Environment

In addition, our method allows to filter tweets by topic, thus it can be applied as starting point to develop a wide variety of analysis regarding other topic or can be used to go deeper into the study of our same topic. Further studies could be focused, for example, on sub-arguments of environment. For example, it would be interesting to filter tweets related to local problematics (i.e. air pollution) rather than to global issues (i.e. global warming) for more detailed longitudinal and spatial studies of the sentiment. This extremely detailed information could be used to study the sentiment on a small scale and, at the same time, to explore how much people care about big themes such as earth health. In this way, we are able to capture the population feelings, to link this to national and/or international policies and events and to identify the main drivers of the inclination and sentiment trends.

Finally, the flexibility of our method can be finalized to create several dictionaries for all the sub-topics connected to a more general phenomenon, such as the well-being (that includes, by nature, different dominion, e.g.: social involvement, health, work status, discrimination; see Toninelli and Cameletti 2018). In this case, selected tweets can be used to study the single dominions and to estimate the subjective well-being and/or how much a single dominion is able to affect the subjective well-being of a population. This will represent an improvement with respect to standard questionnaire-based surveys, such as the European Social Survey<sup>14</sup>. Better, the two types of sources can be integrated. In fact, thanks to the real-time collection of tweets and of similar types of big data, it will be possible to obtain timely information about a multidimensional phenomenon such as the well-being with a very high temporal and spatial resolution. These results can be of high value for evaluating the interventions of policy makers, for measuring the effectiveness of advertising campaigns, for studying a lot of other socio-demographic phenomena.

**Funding:** This work was supported by the University of Bergamo [grant: 60% University Funds, “STaRs - Azione 3: Outgoing Visiting Professor 2019” project].

**Dataset & code:** Datasets used and produced in preparing this paper and R codes are available on Github repository, 2019, link: [https://github.com/silviafabris/Twitter\\_dictionary\\_based\\_classification](https://github.com/silviafabris/Twitter_dictionary_based_classification).

## References

- Ahmed W, Bath P A, Sbaffi L, Demartini G (2019) Novel insights into views towards H1N1 during the 2009 Pandemic: a thematic analysis of Twitter data. *Health Information and Libraries Journal*, 36(1): 60–72
- Alessa A, Faezipour M (2019) Preliminary Flu Outbreak Prediction Using Twitter Posts Classification and Linear Regression With Historical Centers for Disease Control and Prevention Reports: Prediction Framework Study. *JMIR Public Health Surveill*, 5(2), <https://doi.org/10.2196/12383>
- Asghari M, Sierra-Sosa D, Elmaghraby AS (2020) A topic modeling framework for spatio-temporal information management. *Information Processing & Management*, 102340, <https://doi.org/10.1016/j.ipm.2020.102340>
- Baylis P, Obradovich N, Kryvasheyeu Y, Chen H, Coviello L, Moro E, Cebrian M, Fowler JH (2018) Weather impacts expressed sentiment. *PloS one*, 13(4)
- Budiharto W, Meiliana M (2018) Prediction and analysis of Indonesia Presidential election from Twitter using sentiment analysis. *Journal of Big Data* 5(1): 1–10
- [dataset] Schlosser S, Cameletti M, Toninelli D, Tweets datasets used in preparing this paper, Github repository, 2019, link: [https://github.com/silviafabris/Twitter\\_dictionary\\_based\\_classification](https://github.com/silviafabris/Twitter_dictionary_based_classification)

---

<sup>14</sup> For further information, see: <https://www.europeansocialsurvey.org/>.

## Dictionary-based Classification of Tweets About Environment

- Cody E M, Reagan A J, Mitchell L, Dodds P S, Danforth C M (2015) Climate change sentiment on Twitter: An unsolicited public opinion poll. *PLoS ONE*, 10(8): 1–18
- Edo-Osagie O, Smith G, Lake I, Edeghere O, De La Iglesia B (2019) Twitter mining using semi-supervised classification for relevance filtering in syndromic surveillance. *PLOS ONE* 14(7): e0210689. <https://doi.org/10.1371/journal.pone.0210689>
- Foucault N, Courtin A (2016) Automatic Classification of Tweets for Analyzing Communication Behavior of Museums, Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16), 3006-3013, <https://www.aclweb.org/anthology/L16-1480>
- Frenda S, Ghanem B, Montes-y-Gómez M, Rosso P, Pinto D, Singh V (2019) Online Hate Speech against Women: Automatic Identification of Misogyny and Sexism on Twitter. *Journal of Intelligent & Fuzzy Systems*, 36(5): 4743–4752
- Gesualdo F, Stilo G, Agricola E, Gonfiantini MV, Pandolfi E, Velardi P, Tozzi AE (2013) Influenza-Like Illness Surveillance on Twitter through Automated Learning of Naïve Language. *PLoS ONE* 8(12). <https://doi.org/10.1371/journal.pone.0082489>
- Ghafarian SH, Yazdi HS (2020) Identifying crisis-related informative tweets using learning on Distributions. *Information Processing & Management*, 57, 2, 102145
- Grimmer J, Stewart B (2013) Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts. *Political Analysis*, 21(3): 267–297
- Grzeża M, Becker K, Galante R (2020) Drink2Vec: Improving the classification of alcohol-related tweets using distributional semantics and external contextual enrichment. *Information Processing & Management*, 102369, <https://doi.org/10.1016/j.ipm.2020.102369>
- Harb JGD, Ebeling R, Becker K (2020) A framework to analyze the emotional reactions to mass violent events on Twitter and influential factors. *Information Processing & Management*, 102372, <https://doi.org/10.1016/j.ipm.2020.102372>
- Hawkins JB, Brownstein JS, Tuli G, Runels T, Broecker K, Nsoesie EO, McIver DJ, Rozenblum R, Wright A, Bourgeois FT, Greaves F (2016) Measuring patient-perceived quality of care in US hospitals using Twitter. *BMJ Quality & Safety*, 25, 404-413.
- Kozłowski D, Lannelongue E, Saudemont F, Benamara F, Mari A, Moriceau V, Boumadane A (2020) A three-level classification of French tweets in ecological crises. *Information Processing & Management*, 57, 5,
- Kudugunta S, Ferrara E (2018) Deep neural networks for bot detection. *Information Sciences* (467), 312-322. <https://doi.org/10.1016/j.ins.2018.08.019>
- Liu X, Burns AC, Hou Y (2017) An Investigation of Brand-Related User-Generated Content on Twitter. *Journal of Advertising*, 46(2), 236-247. <https://doi.org/10.1080/00913367.2017.1297273>
- Mitchell L, Frank M R, Harris K D, Dodds P S, Danforth C M (2013) The Geography of Happiness: Connecting Twitter Sentiment and Expression, Demographics, and Objective Characteristics of Place. *PLoS ONE*, 8(5)
- Mottaghinia Z, Feizi-Derakhshi M, Farzinvash L, Salehpour P (2020) A review of approaches for topic detection in Twitter, *Journal of Experimental & Theoretical Artificial Intelligence*, <https://doi.org/10.1080/0952813X.2020.1785019>
- Nielsen FÅ (2011) A new ANEW: Evaluation of a word list for sentiment analysis in microblogs. Proceedings of the ESWC2011 Workshop on ‘Making Sense of Microposts’: Big Things Come in Small Packages, Keraklion, Crete, Greece, 93-98
- Nizzoli L, Avvenuti M, Cresci S, Tesconi M (2019) Extremist Propaganda Tweet Classification with Deep Learning in Realistic Scenarios. Proceedings of the 10th International ACM Conference on Web Science (WebSci'19), Boston, USA, 203–204, <https://doi.org/10.1145/3292522.3326050>
- Pamungkas EW, Basile V, Patti V (2020) Misogyny Detection in Twitter: a Multilingual and Cross-Domain Study. *Information Processing & Management*, 102360, <https://doi.org/10.1016/j.ipm.2020.102360>

### Dictionary-based Classification of Tweets About Environment

- Pruss D, Fujinuma Y, Daughton A R, Paul M J, Arnot B, Szafir D A, Boyd-Graber J (2019) Zika discourse in the Americas: A multilingual topic analysis of Twitter. *PLoS ONE*, 14(5): 1–23
- Reyes-Menendez A, Saura J R, Alvarez-Alonso C (2018) Understanding #worldenvironmentday user opinions in twitter: A topic-based sentiment analysis approach. *International Journal of Environmental Research and Public Health*, 15: 2537
- Samuel J, Ali GGMM, Rahman MM, Esawi E, Samuel Y (2020) COVID-19 Public Sentiment Insights and Machine Learning for Tweets Classification. *Information*, 11(6), 314, 1-22. <https://doi.org/10.3390/info11060314>
- Schlosser S, Toninelli D, Fabris S (2019). Looking for Efficient Methods to Collect and Geolocalise Tweets, in *Smart Statistics for Smart Applications – Book of Short Papers SIS2019*. Milan (IT), 18-21 June 2019, 2019, 1057–1062; <https://it.pearson.com/docenti/universita/partnership/sis.html>
- Schlosser S., Toninelli D., Cameletti M. (2020) Comparing Methods to Retrieve Tweets: a Sentiment Approach. *Proceedings of the CARMA 2020 - 3rd International Conference on Advanced Research Methods and Analytics*, 299-306, <http://dx.doi.org/10.4995/CARMA2020.2020.11653>
- Şerban O, Thapen, N, Maginnisa B, Hankina C, Footb V (2019) Real-time processing of social media with SENTINEL: A syndromic surveillance system incorporating deep learning for health classification. *Information Processing and Management*, 56, 1166–1184.
- Snyder LS, Lin Y, Karimzadeh M, Goldwasser D, Ebert DS, (2020) Interactive Learning for Identifying Relevant Tweets to Support Real-time Situational Awareness. *IEEE Transactions on Visualization and Computer Graphics*, 26(1), 558-568. <http://dx.doi.org/10.1109/TVCG.2019.2934614>.
- Stowe K, Anderson J, Palmer M, Palen L, Anderson K (2018) Improving Classification of Twitter Behavior During Hurricane Events. *Proceedings of the Sixth International Workshop on Natural Language Processing for Social Media*, 67–75 (<https://www.aclweb.org/anthology/W18-3512.pdf>)
- Toninelli D, Cameletti M (2018) Is Structural Equation Modelling Able to Predict Well-being?, in A. Abbruzzo, E. Brentari, M. Chiodi and D. Piacentino (Eds), *Book of short Papers SIS 2018, Proceedings of the 49th Scientific Meeting of the Italian Statistical Society*, Palermo (IT), 20–22 June 2018. Pearson, ISBN: 9788891910233 (<https://bit.ly/382HPRA>)