

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

Twitter data analysis to assess the interest of citizens on the impact of marine plastic pollution

Otero^{1*}, P., J. Gago¹, and P. Quintas¹.

¹Centro Oceanográfico de Vigo (IEO, CSIC), Subida a Radio Faro, 50, 36390, Vigo, Spain, Phone Number: 0034986492111, Fax Number: 0034986498626

*Corresponding author: pablo.otero@ieo.es

CRediT author statement

Pablo Otero: Conceptualization, Methodology, Software, Data Curation, Formal Analysis, Visualization, Writing- Original draft preparation. **Jesús Gago**: Conceptualization, Resources, Writing- Reviewing and Editing, Funding acquisition. **Patricia Quintas**: Conceptualization, Validation, Writing- Reviewing and Editing.

1 **Abstract**

2 Few studies have mined social media platforms to assess environmental concerns. In this study, Twitter was
3 scraped to obtain a ~140,000 tweet dataset related specifically to marine plastic pollution. The goal is to
4 understand what kind of users profiles are tweeting and how and when they do it. In addition, topic
5 modelling and graph theory techniques have allowed us to identify main concerns on this topic: i) impact on
6 wildlife, ii) microplastics/water pollution, iii) estimates/reports, iv) legislation/protection, and v)
7 recycling/cleaning initiatives. Results reveal a scarce influence of organizations involved in research and
8 marine environmental awareness, so some guidelines are depicted that could help to adjust their
9 communication plans. This is relevant to engage society through reliable information, change habits and
10 reinforce sustainable behaviour. A visualization tool has been created to analyze the results over time.

11
12 **Keywords:** Twitter, social media, marine litter, plastic pollution, topic modelling, COVID-19

13
14 **Highlights:**

- 15 • Twitter is a valuable tool to analyze the social aspects of marine pollution
- 16 • Topic modelling helped to identify 5 main relevant subtopics
- 17 • COVID-19 pandemic impacted the marine plastic pollution topic on Twitter
- 18 • Low presence of academic or environmental bodies compared to personal opinions
- 19 • An interactive app is released to facilitate further analysis

20
21 **Funding:** This work was co-financed by both the European Regional Development Fund through the
22 CleanAtlantic INTERREG (Atlantic Area) project [grant number EAPA_046/2016] and the Agencia Estatal
23 de Investigación (Ministerio de Ciencia e Innovación, Government of Spain) through ANDROMEDA project
24 [grant number PCI2020-112047].

25
26

1 **1. Introduction**

2 Marine litter is a planetary threat, affecting nearly every marine ecosystem globally (GESAMP, 2015). In
3 particular, plastic constitutes more than 80 per cent of marine litter (European Commission, 2018), and it is
4 estimated that quantities from 4.8 to 12.7 million metric tons per year are entering our seas and oceans
5 (Jambeck et al., 2015). Despite the ambitious commitments currently set by several governments to reduce
6 marine litter, Borrelle et al., (2020) estimated that the annual input may reach up to 53 million metric tons by
7 2030.

8 The impacts of plastics on marine ecosystems are broad, including habitat degradation and a wide range of
9 negative effects on marine organisms. The impacts include from wildlife dead due to ingestion, starvation or
10 entanglement in marine litter (Gall and Thompson, 2015), to attaching and drifting invasive species and
11 pathogens (i.e., hitchhiking), among others. The socioeconomic effects are also evident in sectors such as
12 fisheries (e.g., damaged gear during trawling activities or reduction of catches), on tourism due to the
13 presence of beach litter or the economy of coastal areas due to clean up actions (GESAMP, 2015). Indirect
14 effects on human health are still being discussed, including the sources and transport dynamics of antibiotic
15 resistance (Bank et al., 2020) and the still unknown effects of microplastics (<5 mm) along the food chain
16 (GESAMP, 2020). Those microplastics can enter the marine environment as primary microplastics (e.g.
17 manufactured pellets and microbeads) or secondary microplastics after the fragmentation and degradation of
18 larger plastics. The presence of microplastics in all marine environments, including marine biota, has been
19 reported in several scientific studies (Filgueiras et al., 2020; Gago et al., 2020; GESAMP, 2020, 2015).

20 Nowadays, it is impossible even to try imagining our world without plastics, given the extreme importance
21 and the number of functions it has in a broad range of aspects of the industry and everyone's daily life. There
22 are no other manufactured materials whose production has grown as plastic has over the last 70 years (Geyer
23 et al., 2017). From 1950 until 2015, 8,300 million metric tons of plastics have been produced: 30% of
24 products are in use, 10% has been incinerated, only 7% has been recycled and 55% has been discarded
25 (Geyer et al., 2017). It is clear that there is an excess of consumption and that many single use articles could
26 be substituted by other materials. Areas of high population density, poor waste management or lack of
27 environmental education, become factors that favour littering of the aquatic environment by plastics (Duckett
28 et al., 2015; Napper and Thompson, 2020).

29 As a consequence, the potential solutions to mitigate the problem are widespread, and the governance
30 solutions become complex. Government and legislative initiatives, changes in the industry and a greater
31 environmental awareness of citizens are factors that can help reduce the arrival of plastics into the sea (Vince
32 and Hardesty, 2017). As part of this strategy, understanding public perceptions, opinions and knowledge
33 about marine plastic litter issue is a critical step in effectively engaging society and changing human

1 behaviour (Forleo and Romagnoli, 2021).

2 In the last decade, the information about marine litter has circulated from the scientific community to the
3 public, through reports, awareness campaigns, events and informative material of all kinds. The disclosure
4 has contributed to raising a critical conscience in society (Heidbreder et al., 2019; Mitrano and Wohlleben,
5 2020; Vince and Hardesty, 2017). This information has been increasingly echoed in part by generalist media.
6 However, another part of the success must be attributed to Social Media (SM), which has begun to open the
7 eyes of many people regarding some of these environmental threats. SM users have passed the 3.8 billion
8 mark and were estimated that more than half of the world's total population was using SM by mid-2020
9 (Kemp, 2020). Viral messages, photos and videos can reach audiences of millions (Parton et al., 2019), so
10 data created and shared by users on SM platforms have emerged as a potentially useful source of information
11 in marine environmental research, management and conservation (see e.g, Abreo et al., 2019; Becken et al.,
12 2017; Ghermandi et al., 2020; Parton et al., 2019; Retka et al., 2019; Ruiz-Frau et al., 2020).

13 Twitter is one of the most popular SM and microblogging sites with more than 330 million monthly active
14 users worldwide (Kemp, 2020), who post ~500 million comments (the so-called *tweets*) per day with up to
15 280-characters containing their thoughts and opinions. Despite the limited number of characters, the
16 possibility of including links allows to increase the information and reach a higher impact. These tweets are
17 affected by both real-world events and the trends of other messages posted in SM (Zubiaga and Ji, 2014).
18 Nowadays, Twitter is the most important SM on science dissemination where journalists, science
19 dissemination professionals, scientists and many research institutions interact, talk and share science with
20 other colleagues and the public (Collins et al., 2016; Letierce et al., 2010; Mohammadi et al., 2018; Van
21 Noorden, 2014). About 50% of scientists use Twitter to follow conversations or debates about their discipline
22 and about 40% consider this network as a tool to talk about their progress or that of other colleagues (Van
23 Noorden, 2014). The most common perceived benefits of Twitter were the size and diversity of the audience,
24 the ability to network with other scientists and the ability to engage with the public (Collins et al., 2016;
25 Smith, 2015).

26 Twitter has increasingly become a world-wide choice to raise awareness and disseminate information on a
27 variety of topics, as the promotion of cancer screening and early diagnosis through specific campaigns
28 (Plackett et al., 2020; Teoh et al., 2018; Vraga et al., 2018; Yoosefi Nejad et al., 2019), identify cancer
29 barriers and policy solutions (Shimkhada et al., 2021), identify mental health discourses (Budenz et al., 2020;
30 Makita et al., 2021), detect and predict the epidemic of diseases (Dang et al., 2018), analyze pro- and anti-
31 vaccination discourses (Milani et al., 2020), aware about emerging technologies (Li et al., 2017),
32 emergencies (Barker and Macleod, 2019; Martínez-Rojas et al., 2018; Zhou et al., 2021), and so on. Karami
33 et al. (2020) found 38 different topics in more than 18,000 Twitter-related papers published between 2006

1 and 2019, using analysis techniques like sentiment analysis, topic modelling or graph mining, among others.
2 These techniques have been applied here to a dataset of more than 140,000 tweets to analyze the interests of
3 citizens about marine litter. Text mining and natural language processing were used to learn about what
4 people commented on this particular SM and how they did it, identifying dominant topics and analysing the
5 word and hashtag frequencies. Sentiment analysis was also employed to explore which were their feelings,
6 whereas geo-tagged tweets and information from the user profile were combined to know the main hot spots
7 of discussion. Additionally, manual and automatic identification of image content in tweets was conducted.

8 This study aims to describe the interest and awareness of marine pollution by plastics and microplastics in
9 Twitter, to understand the spatio-temporal trends and sentiment of tweets, to distinguish different subtopics
10 from the general discourse. Additionally, those images associated with tweets have been explored to assess,
11 among other things, their suitability for discerning amounts and types of litter, particularly in coastal areas.

12 Taking into account that people engage with information posted by people they trust (Huber et al., 2019;
13 Media Insight Project, 2017), this study will provide new insights to governmental, academia and NGOs
14 involved in marine environmental protection to reanalyze their communication strategy on Twitter.
15 Understanding who tweets about the marine litter issue and how they do it will help institutions to design
16 effective communication on this channel to reinforce the commitment of users who are already engaged,
17 facilitate greater public understanding of solutions and enable action.

18 **2. Methodology**

19 This study aims to perform an exploratory analysis of a collection of tweets that cover frequency analysis,
20 sentiment analysis, graph theory and topic modelling, among others. Scraping and data mining techniques
21 involve different steps from data acquisition and data cleaning to data analysis (see Figure 1). In addition, an
22 automatic classification image analysis technique has been tested with the aim of characterizing litter in
23 coastal areas.

24 **2.1. Dataset creation**

25 Figure 1 shows a flow chart of the methodology applied during dataset creation and data analysis phases. The
26 first step consists of collect data from Twitter database. The set of streaming APIs offered by Twitter gives
27 developers low latency access to Twitter's global data, which include the tweet text along with the associated
28 metadata (post time, geographical coordinates if geolocation is enabled, information about the user profile,
29 etc.). In this study, free Twitter's standard search API v1.1 (search/tweets) was used for simple queries
30 against the indices of recent or popular tweets and behaves similarly to, but not exactly like the search UI
31 feature available in Twitter mobile or web clients. The Twitter Search API works as a keyword search

1 method against a sampling of recent tweets published in the past 7 days (further details available in
2 <https://developer.twitter.com/en/docs/twitter-api/v1/tweets/search/api-reference/get-search-tweets>); words
3 and not hashtags were used to perform a query to capture a higher number of tweets. To analyze a longer
4 period, a Python script scrapped the Twitter service every week.

5 Several queries were done to retrieve original tweets (not retweets) by a combination of the keywords
6 ‘*plastic*’ or ‘*microplastic*’ with at least one keyword related to the marine environment: ‘*ocean[s]*’, ‘*sea*’,
7 ‘*beach*’, ‘*coast*’ and/or ‘*marine*’. To ensure representative results of the global use made of this platform,
8 main languages used on Twitter were exposed to queries: English, Japanese, Spanish, Portuguese, French,
9 Italian, Malaysian, German, Turkish, Thai, Korean and Indi (sorted in descending order in our dataset).
10 Tweets were retrieved weekly during ~8 months (Mar 19 – Nov 16, 2020). Table 1 displays a few tweets
11 from the final dataset of 147,552 tweets. The dataset needs to be cleaned up before analysis. Tweets usually
12 contain colloquial and informal sentences, URLs, emojis and emoticons. Therefore, some cleaning is needed
13 to facilitate its understanding and analysis. At the same time that the original text is preserved for each tweet
14 for future reference, a “*sanitized*” text is added to the dataset. Hashtags, user mentions, URLs, media and
15 symbols are stripped out from the full text. Although only original tweets were retrieved, preventive cleaning
16 of “RT” (retweet) at the beginning of the text and symbols such as *at signs* were stripped out from the full
17 text. The remaining text was lower cased and automatic translation of those non-English tweets was done
18 using the TextBlob library in Python (<https://textblob.readthedocs.io/>). Internally, TextBlob relies on Google
19 Translate's API. To skip the rate limit, a delay was done between consecutive queries to the API. Finally,
20 English text was cleaned of grammatical contractions (e.g., “*ain't*” to “*is not*”, “*I'll*” to “*I will*”); with this
21 purpose, a set of 125 contractions was used (<https://github.com/PabloOtero/twitter-python>). While some
22 abbreviations and acronyms may be common across all SM sites, others are unique to the microblogging
23 platform. Some terms of this Twitter lingo were also fixed (e.g., “*u*” to “*you*”, “*ya*” to “*yeah*”).

24 With clean text, it is easier to apply sentiment analysis and store the results along with our dataset. The
25 algorithms of sentiment analysis mostly focus on defining opinions, attitudes, and even emoticons in a
26 corpus of texts. In this study, the approach of the TextBlob package based on a pre-defined set of categorized
27 words was used. The sentiment property returns both *polarity* and *subjectivity*. The polarity score is a float
28 within the range [-1.0, 1.0]. The subjectivity is a floating value within the range [0.0, 1.0] where 0.0 is very
29 objective and 1.0 is very subjective. Two sentiment scores were computed for each tweet: one based only on
30 words and the other one on words plus emojis and emoticons; the meaning of an emoji will depend on the
31 context of the current text. Whereas emoticons are handled well by TextBlob, emojis were searched in the
32 message based on a Unicode list (<http://www.unicode.org/emoji/charts/full-emoji-list.html>) and converted to
33 official name or any known short name (<https://github.com/alexmick/emoji-data-python>).

1 Whenever possible, the spatial information associated with the tweet was also obtained and added to the
2 dataset. Geographical coordinates were extracted from tweet metadata (geotagged tweets) in those cases in
3 which the user had the geolocation enabled in its device, something that only happened in 1.98% of the
4 cases. Coordinates were exposed to Nominatim API (<https://nominatim.org/>) to perform reverse geocoding
5 on OpenStreetMap (<https://www.openstreetmap.org/>) and obtain city, state —when possible— and country
6 data fields. From those non geotagged tweets, location information from the user profile was obtained and
7 forward and reverse geocoding performed through the API. Here, it is necessary to understand that the user
8 does not necessarily have to indicate an existing or recognizable place. In fact, although 70.6% of the tweets
9 presented a location in the user's profile, only 56.7% had a valid position on the globe.

10 As summarized in Figure 1, further cleaning of the dataset may be necessary to facilitate some analysis and
11 the interpretation of results. This is the case of frequency analyses where the social bots influence the volume
12 data. Although there is no universally agreed-upon definition of a bot, they can be considered malicious
13 actors that create inauthentic social media accounts partially controlled by algorithms. Most automated bots
14 reply or post tweets or simply follow other users based on triggers or according to some scripted patterns
15 (e.g., retweeting all messages from certain accounts). To detect social bots, the popular API Botometer® v4
16 (Sayyadiharikandeh et al., 2020) was used. Botometer is a machine learning algorithm trained to calculate a
17 score where low scores indicate likely human accounts and high scores indicate likely bot accounts. From
18 several types of scores provided by Botometer, we used the "overall score" based on a comparison of several
19 models trained on different kinds of bots and human accounts and language-independent (“universal”). In the
20 [0-1] range, here we considered 0.6 as a limit to consider an account as a bot.

21 In the same way that the detection and elimination of bots are important, it is also necessary to discard as
22 much as possible those tweets that are not directly related to marine plastic pollution. As a notable example,
23 are those tweets that refer to the British virtual band “Gorillaz” —with about 1 million followers on
24 Twitter— and their studio album “*Plastic Beach*”. Thanks to the Latent Dirichlet Allocation (LDA)
25 technique (Blei et al., 2003) that will be described later, we could determine the most frequent words used
26 when talking about this subtopic (e.g., ‘song’, ‘album[s]’, ‘music’, ‘demon’, ‘humanz’, ‘listening’) and
27 consequently, eliminate the related tweets.

28 **2.2. Data analysis**

29 **2.2.1. Topic modelling**

30 Topic modelling is an efficient and systematic approach to analyze thousands of documents in a few minutes
31 (Karami et al., 2020). Among topic models, LDA is a valid and widely used generative probabilistic model
32 (see e.g., Blei et al., 2003). LDA identifies semantically related words, which occur together in multiple

1 documents (i.e., tweets) of a corpus (i.e., our preprocessed dataset as shown in Section 2.1). As a result,
2 several groups of multinomial distributions over the terms in the vocabulary of the corpus represent the
3 topics. To interpret a topic by human intuition as a meaningful “*theme*”, one typically examines the top terms
4 in a ranked list of the most probable terms in that topic. The problem with this method is that common terms
5 in the corpus often appear near the top of several lists, making interpretation difficult. For this reason, these
6 lists were sorted by relevance according to Sievert and Shirley (2014). They defined the relevance r of word
7 ω to topic k given a weigh parameter λ [0-1] as:

$$8 \quad r(\omega, k|\lambda) = \lambda \log(\phi_{k\omega}) + (1 - \lambda) \log\left(\frac{\phi_{k\omega}}{p_\omega}\right) \text{ [Equation 1]}$$

9 where $\phi_{k\omega}$ denotes the probability of term ω for topic k and p_ω the marginal probability of term ω in the
10 corpus. Whereas $\lambda = 1$ shows the classical ranking of terms by their probability in the topic, a lower λ also
11 weights by the probability of appearing in the corpus. At the other extreme, $\lambda = 0$ classifies terms by the ratio
12 in a logarithmic scale of their probability within a topic to its marginal probability across the corpus, also
13 known as *lift* (Taddy, 2011); this is, ranks words that appear exclusively in that topic but not in the others.
14 Thus, playing to vary λ can help to better define the associated topic.

15 In this study, the technique was only applied to English words after removing from the corpus those
16 keywords used in data acquisition; in this way, the subtopics produced by this technique are meaningful and
17 not dominated by the same keywords. Stop words, which are the most common without significant
18 contextual meaning in a sentence (e.g., “*a*”, “*the*”, “*and*”, “*but*”, and so on) were also filtered out from the
19 corpus. Those ampersands (&) written in tweets via a mobile device appear in the document as “*amp*” with
20 no inherent meaning, so they were removed as well. Apostrophes were also deleted and words in plural were
21 converted to the singular as far as possible.

22 The number of topics must be chosen before LDA is run; however, it is unclear how many topics the dataset
23 should be divided into. A low number of topics can cause the loss of a detailed view of the text when
24 merging topics. Alternatively, a high number of topics can lead to too many top words being shared and
25 make interpretation difficult. In this study, different tests were run with some topics varying from 5 to 15 and
26 finally, 10 words in 6 topics were used to train the model. To better understand the underlying fitted LDA
27 model, the LDAvis tool (Sievert and Shirley, 2014) over Python was used
28 (<https://github.com/bmabey/pyLDAvis>). This tool allows flexibility in exploring topic-term relationships
29 using *relevance*.

30 **2.2.2. Graph theory analysis**

31 To complement topic modelling, analysis of networks using graph theory was performed. A graph (network)

1 is a collection of vertices (nodes) with a collection of edges that are connections between the different
2 vertices in a network. In this study, nodes are represented by words, while edges illustrate the connections
3 between words in the same tweet and the frequency of those connections. Within the network, it is possible
4 to distinguish communities. A community is defined as a group of nodes where the density of the edges
5 between the nodes inside the group is greater than the connections with the rest of the network. To find
6 communities in our network, a semi-synchronous label propagation method was used (Cordasco and
7 Gargano, 2011). This method combines the advantages of both synchronous and asynchronous models. If a
8 score is given to the number of links between two nodes and the process is repeated for the complete network
9 landscape, the *modularity* —a measure of the strength of the division of a network into communities— can
10 be computed. Networks with high modularity have dense connections between the nodes within communities
11 but sparse connections between nodes in different communities. According to Clauset et al. (2004),
12 modularity can be computed as:

$$Q = \sum_{c=1}^n \left[\frac{L_c}{m} - \left(\frac{k_c}{2m} \right)^2 \right] \text{ [Equation 2]}$$

14 where the sum iterates over all community c , m is the number of edges, L_c is the number of intra-community
15 links for community c and k_c is the sum of degrees of the nodes in the community c . Modularity ranges from
16 -1 to 1, and the higher the value, the better the community structure.

17 On the other hand, different types of centrality measures can be used to identify which nodes are the biggest
18 influencers on the network. Here, an eigenvector centrality, which is based on the centrality of its neighbours
19 was used (Newman, 2010). A node with a high score will influence multiple nodes, which in turn are highly
20 connected. The advantage of this method is that it can highlight nodes that exercise control behind the
21 scenes. For the creation, manipulation, and study of the structure, dynamics, and function of our network the
22 library NetworkX was used (<https://github.com/networkx/networkx>; Hagberg et al., 2008). By default, the
23 layout of the nodes and edges is automatically determined by the Fruchterman-Reingold force-directed
24 algorithm (Fruchterman and Reingold, 1991). Frequencies of word pairs, also called bigrams, were analyzed.

25 **2.2.3. Image analysis**

26 Photos associated with tweets were also downloaded and stored for image content analysis (6,172 images
27 only in English tweets). Duplicated images were removed after comparing their associated Message Digest
28 Algorithm 5 (MD5) hash values. The Computer Vision API (v3.1) from Azure Cognitive Services
29 (<https://azure.microsoft.com/es-es/services/cognitive-services/>) was used to process and obtain information
30 from images. This free tool allows, among other features, to estimate the dominant and accent colours,
31 categorize the content of the images, tag and create a short description. In the present study, the goal was to
32 use this tool to filter images based on categories and tags to assess the type of media content uploaded by

1 users, as a previous step to object detection techniques in future studies.

2 **3. Results and Discussion**

3 We first present results from a social perspective to know, among other aspects, which languages were used
4 the most, from where people were tweeting, when they were active, the most frequent words and hashtags,
5 the main subtopics, the positive or negative feelings and who capitalized on the conversation. Second, image
6 analysis was performed to determine the type of content most used and to assess whether this information
7 could be used to monitor coastal areas directly from images present on Twitter.

8 **3.1. Regional analysis**

9 Figure 2 displays a heatmap for the dataset where positions of each tweet were obtained either from the
10 metadata in the case of geotagged tweets or estimated from the user's profile. The spatial information
11 obtained in this way accounted for 56.7% of the total volume of tweets. The map shows a greater
12 concentration of tweets on the coasts of the USA, Japan, Western Europe, the west coast of South America,
13 Indonesia and the west coast of Australia. Hot spots in the map can be compared with the number of tweets
14 per country shown in Figure 3a. The top four countries were the USA with 16,111 tweets, the UK with 9,908,
15 Japan with 5,909 and Canada with 3,548 tweets. These countries accounted for 52.4% of the total tweets
16 with associated spatial information in their metadata.

17 To achieve a vision as global as possible we have made a multilingual approach, in contrast with the majority
18 of published studies that queried Twitter with hashtags or only with keywords in English. The frequency of
19 languages in our dataset was English (63.1%), Japanese (16.5%), Spanish (8.5%), Portuguese (2.9%), French
20 (2.8%), Italian (1.9%), Malaysian (1.8%), German (0.8%), Turkish (0.6%), Thai (0.5%), Korean (0.4%) and
21 Hindi (0.2%), a list that does not coincide in order with the classification by most widespread native languages
22 (Eberhard et al., 2019) nor the usage statistics of content languages for websites (W3Techs, 2020). If we
23 narrow the list in Figure 3a to only English-speaking countries, the top 4 countries by volume of tweets
24 become the USA, UK, Canada and Australia, the same countries with a greater number of comments in
25 Twitter on the climate change issue (Dahal et al., 2019).

26 To compare with the population size, the tweet volume was normalized by the population of each country
27 and by the corresponding maximum ratio from these top-25 countries (Figure 3b); the list is restricted to
28 prevent a high bias by little populated countries with a large relative volume of tweets (e.g., small island
29 states). This figure reveals that the UK was the country with the highest number of tweets per capita
30 followed by Ireland, Canada, New Zealand and Spain. Although Japanese was the second language in our
31 dataset, Spain appeared as the first non-English-speaking country with the highest relative weight. A third

1 approach could be taken to weigh based on the digitization of the country as well as the engagement of this
2 social network. Here, the tweet volume per country was divided by the number of active users on Twitter.
3 Although we have only had access to data from countries with the highest number of active users
4 (STATISTA, 2020), it was enough to verify that the order of the previous lists would be modified, becoming
5 now the top 4 countries: UK (with 16.6 Million Active Users, MAU), Spain (7.5 MAU), France (7.9 MAU)
6 and Germany (5.4 MAU). The USA occupied the fifth position as the interest in this topic dissolved among
7 its high number of active users (68.7 MAU). Saudi Arabia occupied the eighth position in the number of
8 active users however it was not represented in our dataset, so Arabian should be also considered in the
9 queries to Twitter in future studies.

10 The difference in data volume between countries is a combination of population size, the degree of
11 digitization, the number of active users on this social network and, finally, interest in this specific topic. For
12 example, India is the third country in the world in terms of active users (18.9 MAU) on Twitter surpassing
13 slightly UK. If we focus on a highly topical issue during the data acquisition period such as COVID-19,
14 India was ranked third globally in line with its number of active users (Banda et al., 2020). However, India
15 occupied the ninth position in data volume and was the thirtieth in terms of tweets per capita in our study,
16 contrasting with the UK that demonstrated to be the country with the greatest interest in the marine plastic
17 pollution topic on Twitter. This lower relative interest in comparison with its number of active users was also
18 observed in countries such as Japan or Brazil, among others.

19 Over a quarter of tweets with spatial reference in the dataset came from the USA which invited to deepen the
20 analysis in this country. California was the state with the largest volume of domestic tweets (15.7%) followed
21 by New York (9.7%), Kansas (8.1%) and Florida (7.8%). Dahal et al., (2019) found that the northeast region
22 had a relatively high amount of climate change discussion and this could be caused, among other hypotheses,
23 by the cultural and political differences, since climate change is treated as a political issue by many Twitter
24 users.

25 A high amount of tweets in densely populated coastal states was one of the expected results in this study, yet
26 the inland state of Kansas was surprising. In this state, there was not high activity of any particular user or
27 group of users. Nor had a higher immersion been observed in environmental campaigns held online.
28 Therefore, it will be interesting to analyze the domestic behaviour of USA in future studies that include a
29 longer period.

30 **3.2. Temporal analysis**

31 The temporal evolution is impacted by different events as depicted in Figure 4. The largest peak corresponds
32 to the celebration of the World Environment Day on 5th June 2020 with the lesser impact of other world-wide

1 celebrations and campaigns with hashtags such as #EarthDay (22nd June), #WorldSeaTurtleDay (16th June) or
2 #PlasticFreeJuly (1st-3rd July). The second and third largest peaks on the series are related to comments on
3 environmental reports published by The Pew Charitable Trusts and SYSTEMIQ (2020) and OCEANA
4 (Warner et al., 2020) on 14th July and 19th November, respectively. Scientific publications in high impact
5 journals also impacted the volume of tweets, like the study of Pabortsava and Lampitt (2020) on 18th August,
6 Borrelle et al. (2020) on 6th October and Law et al. (2020) on 5th November. The comments on Twitter about
7 the scientific studies echo news in high-audience media (e.g. *The New York Times*, *FOX News*, *The*
8 *Guardian*, *The Economist*, etc.), which explains a certain delay between the publication date on scientific
9 journals and their peak in the time series on Twitter.

10 The acquisition data began with an exceptional situation due to the COVID-19 pandemic, with national
11 lockdown measures, particularly in most of Europe, Asia and South America, followed by a period of greater
12 freedom over the summer months. Twitter is what's happening and what people are talking about right now.
13 In this sense, the temporal analyses were affected by the confinement period with people at home and face to
14 face activities focused on marine litter cancelled (e.g. beach clean-ups, events, etc.).

15 Nevertheless, the conversation around this topic remained active on Twitter, probably due to the concern
16 about the increase of single-use plastics (especially masks and gloves) during the COVID-19. Thus, before
17 1st June, the average number of daily tweets was 351.5+/-104.6, whereas after this date the average was
18 significantly higher 472.7+/-172.69 ($p<0.05$). The COVID-19 situation also pulls the comments, with a peak
19 after the article in "*The Economist*" entitled "*COVID-19 has led to a pandemic of plastic pollution*" and
20 published on 22nd June. (<https://www.economist.com/international/2020/06/22/covid-19-has-led-to-a-pandemic-of-plastic-pollution>).
21

22 If tweets are grouped in time hour slots and days of the week (see Figure 5), results show the highest number
23 of posts on business days, between 12 and 18 UTC, and increases as the workweek progress. As expected,
24 there are differences if the analysis is by time zone, mainly due to differences in social habits and work
25 schedules. Peaks often coincide with catching up on the coffee break, lunchtime, and the time people are on
26 their way home. For example, in Japan (figures not shown), the activity is high from 7 am to midnight and
27 during all days of the weeks, with some peaks at noon and 6 pm. In Spain, Twitter activity in this topic is
28 high from 9 am to 9 pm, with peaks at noon and 2 pm; the increase in activity on Sunday is noteworthy. In
29 the USA, the activity is high between 7 am and 6 pm, mainly during business days and particularly on
30 Thursdays. In this country, there are also important differences between the Eastern and Central Time Zones.
31 Habits and work breaks determine the time of use of this SM and being aware of this reality is relevant to
32 increase the engagement, particularly taking into account the "*short lifespan*" of a tweet (Wilson, 2019).

33 3.3. Sentiment analysis

1 Figure 4 also shows the volume of positive and negative tweets per day. Positive tweets (n=67,470) always
2 exceeded the negative ones (n=33,612). An increase of positive tweets was noticeable during the celebration
3 of both the Earth Day and World Environment Day. In contrast, the volume of negative tweets in social
4 networks increased during the celebration of Sea Turtle Day with many references to the entanglement of
5 turtles. Negative tweets were also accompanied by references to Pew's report (The Pew Charitable Trusts
6 and SYSTEMIQ, 2020), which warn about poor environmental conditions.

7
8 Figure 6a shows a histogram of the polarity of tweets at 0.25 bin intervals. On average, tweets exhibit
9 significantly greater sentiment when emoticons are included (0.085 ± 0.264) than when these are not taken
10 into account (0.065 ± 0.256) (paired t-test; $p < 0.01$). Negative tweets are significantly more objective
11 (0.435 ± 0.265) than positive tweets (0.516 ± 0.202) (2-sample t-test; $p < 0.01$), although the frequency
12 distribution encourages us to interpret this result with caution (see Figure 6b and 6c).

13
14 The comparison of average sentiment between the 10-top countries (those with more than 2,000 tweets) also
15 shows significant differences (ANOVA; $p < 0.01$), with Spain the most positive country (0.123 ± 0.262) and
16 Japan the less positive (0.038 ± 0.219). Internally in the USA, no significant differences were found on
17 average sentiment between the west and east coasts of the USA and neither between the coastal and inland
18 states.

19 20 **3.4. User's activity and engagement**

21
22 Just as important as knowing where and when people tweet about marine plastic pollution, it to know who
23 does it, how often, and the success of their message. From 81,664 users that composed our dataset, we have
24 sorted the top 100 users by the number of tweets, by engagement and by the number of followers. To avoid
25 ethical and privacy problems, the data showed here (Figure 7) has been aggregated in different categories;
26 explicit mentions in the text are only related to large organizations or companies and not individuals.

27 The first way to categorize relies on a binary classification between bots and human-like user profiles. Based
28 on results from the Botometer API, those user profiles with an overall bot score greater than 0.6 and
29 probability above 0.8 were directly classified as bots and the rest were supervised based on the number of
30 published tweets, number of followers, the type of content, etc.; this binary classification can sometimes be
31 subjective, as bots are becoming more and more refined. In turn, human-like profiles were classified as
32 companies, individuals, NGOs/foundations/nonprofit organizations, academic institutions, official organisms
33 or initiatives/projects.

34 As expected, bots are the main group when users are classified by tweet volume (32%), followed by

1 individuals (28%) that double the rest of the categories. Generally, these individual users are not celebrities
2 or influencers, in contrast to the list of users with more engaging tweets like actors, musicians, soccer
3 players, writers, politicians or even astronauts. Companies are the second category in terms of successful
4 tweets, particularly those coming from big media (BBC, CNN, New York Times, The Economist, ABC,
5 Globo News, Huff Post, Le Monde, The Guardian, etc.); nothing surprising considering that they have the
6 highest number of followers. The Ocean Clean Up, Oceana, Earth Day Network, WWF Japan, No Plastic
7 Waste or Greenpeace are some examples of NGOs/Foundations/Nonprofit organizations with tweets with
8 large engagement. Within the initiatives/projects category, Lost at sea, Aplastic Planet and Blue Planet
9 Society are some examples. The presence of both academic institutions and official bodies responsible for
10 the care and protection of the environment is scarce.

11 **3.5. Hashtags and topic modeling**

12 The use of hashtags has the advantage of classifying tweets within a certain topic with some independence of
13 the language. Figure 8 shows the most used hashtags in the complete dataset. Results show a classification of
14 events like the World Ocean Day or Plastic Free July, being the most used hashtag *#plastic* followed by
15 *#plasticpollution*. The positive sentiment outweighs the negative in all hashtags, although neutrality
16 predominates. Although hashtags allow immediate thematic classification of the tweet, they may not be
17 sufficient to determine subtopics. Here, the corpus of the dataset must be used. To simply analyze the
18 occurrence of words, only English tweets were processed to avoid problems related to automatic translation,
19 which could alter the meaning of the word or use different synonyms in the translation. The words *pollution*
20 and *waste* appear in similar frequency (9,121 and 9,115 occurrences, respectively), followed by *use* (8,514),
21 *help* (7,477) and *people* (6,781).

22 Table 2 shows the top-10 most probable words in the 6 topics generated by LDA, with a weigh parameter λ
23 of 1 and 0.4. Whereas $\lambda = 1$ shows the ranking of terms by their probability in the topic, a lower λ also
24 weights by the probability of appearing in the corpus. At the other extreme, $\lambda = 0$ ranks words that appear
25 exclusively in that topic but not in the others. Thus, decreasing the value of λ means that less frequent and
26 more exclusive words of that topic will rise in the ranking, although that does not necessarily imply that it
27 helps to better define the topic by a human. From model results, the following six topics ordered from larger
28 to lower marginal topic distribution were defined:

- 29 i) “*Impact on wildlife*” (20.4% of tokens) that defines concerns about the impact of marine litter, especially
30 plastic bags and straws on marine biota, with particular mention to turtles;
- 31 ii) “*Microplastics/Water pollution*” (20.2%) referring to the pollution produced by plastics and microplastics
32 derived from items such as plastic bottles and bags;
- 33 iii) “*Estimate of quantities/Reports*” (17.5%) with comments on the amounts of marine litter that impact the

- 1 environment based on news from recent studies and reports;
- 2 iv) “*Legislation/Protection*” (14.6%) concerning problems of legislation and the need for global treaties to
3 tackle this problem;
- 4 v) “*Recycling/Cleaning initiatives*” (11.1%) with comments on citizen initiatives, private companies and
5 NGOs related to the collection and cleaning of marine litter and the reduction of single-use plastics, among
6 others;
- 7 vi) tweets with comments about the album “*Plastic beach*” of the British virtual music band “*Gorillaz*”
8 (16.2%), completely unrelated to the subject of this study.

9 Table 2 also includes a selection of words that, without being in the highest positions in the ranking, may be
10 relevant. Thus, for example, words like *birds*, *mammals* or *entangled* also define the topic “Impact on
11 wildlife”, whereas words like *butts*, *cigarettes* or *fibers* complement the definition of the topic
12 “*Microplastics/Water pollution*”. Words about the COVID-19 pandemic situation were mainly related to the
13 topic “*Estimate of quantities/Reports*”.

14 According to Mehrotra et al. (2013), the performance of the topic models produced by LDA on Twitter data
15 is significantly improved when tweets are aggregated by some common factor to produce pseudo-documents
16 for the corpus. Thus, we have also merged documents (tweets) from the same users before performing LDA,
17 similar to a recent study by Dahal et al. (2019). Same hyperparameters and corpus cleaning methods were
18 applied to both methods. To compare the quality by topic between both models, the metric of UMass
19 coherence (Röder et al., 2015) was examined.

20 The UMass coherence measure assesses topic quality by looking at how frequently words within a topic co-
21 occur in the corpus. The average UMass coherence of the author-pooled LDA was -2.96 and the classical
22 LDA was -4.63. Despite the author-pooled LDA performed better in statistical terms, we found easier to
23 interpret topics determined by classical LDA. For example, The LDAvis tool (Sievert and Shirley, 2014)
24 used to interpret results showed overlap of token clusters from three topics in the author-pooled LDA and
25 some of the words assigned to the topics had relatively little meaning. LDA is fundamentally a statistically
26 trained model and its performance does not always directly translate to better human interpretability. In fact,
27 topics in Table 2 are indeed meaningful, which was the desired result. As Twitter users add hashtags to align
28 their tweets with a specific topic, it is expected that pooling tweets that share the same hashtag would
29 produce better topic models (e.g., Dahal et al., 2019; Steinskog et al., 2017). However, only a fraction of
30 tweets contain hashtags (28%) and the selection of a particular one or group of them would imply abandonee
31 the most of our dataset. Pooling by hashtags is however straightforward in the analysis of other global and
32 well identified issues, as could be the #MeToo movement (e.g., Goel and Sharma, 2020; Manikonda et al.,
33 2018) or #COVID-19 (e.g., Xue et al., 2020).

1 In conjunction with topic modelling, network visualization has been used to make sense and explore our
2 dataset. Figure 9 shows networks of the 100 most frequent pairs of co-occurring words (bigrams) in tweets
3 on marine pollution by plastics from March 19th to June 1st, 2020, when a large part of countries world-wide
4 were in strict confinement measures due to COVID-19 pandemic. This period is particularly interesting
5 because without differing too much from the network graph for the entire study period (not shown), it
6 highlights relationships with words typical of the pandemic status: *lockdown*, *COVID-19* and *pandemic*. The
7 label propagation method detects 14 communities in the subset, with a modularity value of 0.39. This poor
8 value can be partially explained by the shortlist of pairs used to build the network. Many of these
9 communities are composed of only a pair of words and the largest by 25 nodes. The word *plastic*, our main
10 keyword in this study, is crucial to the network and its centrality value is used to normalize the rest of the
11 nodes. Its community is composed of words such as *bag*, *bottle*, *debris*, *pollution* or *straws*, but also the word
12 *pandemic* belongs to this community. This last word links to a different community formed by the words *led*
13 and *COVID-19*. The *ocean* is the second word in the network with the highest eigenvector value, meaning
14 that it is highly influencing other strongly-connected words in the net. As an example of other useful
15 information that can be extracted from the graph, *plastic* links through *pollution* with a community related to
16 *legislation* and the need of *treaties* at a *global* scale; the low centrality values of this community suggest that
17 it is a well-defined community separated from the rest.

18 Finally, filtering our dataset by the presence of the word *stomach* in the corpus could help to list the marine
19 species in which plastic ingestion has been observed; alerting even before there is a scientific publication
20 with the observation. Although this word was not among the most frequent ones, 506 tweets were found
21 related to this subtopic.

22 **3.6. Image analysis**

23 A total of 33,285 tweets (24%) contained associated media resources. It is well known that including an
24 image in a tweet increases engagement (Wadhwa et al., 2017). Publishing a tweet with visual content makes
25 the publication more attractive and suggestive for the user. The image catches the interest of the user, and
26 surprises and acts as a claim. The message takes up more physical space on the user's screen and helps the
27 message to be better understood (Polinario, 2016). Tweets with pictures generate greater engagement
28 independently of their content (Carrasco-Polaino et al., 2019).

29 To know which tweets with associated images aroused the most interest among users, the tweets were
30 classified by the total number of retweets and favorites. In our dataset only the original tweets were kept and
31 not the retweets, therefore, only direct interactions with the original tweet were taken into account. A top-100
32 list was done and engagement examined. Focusing on the top 10 tweets on this list, 2 of them belonged to
33 accounts that were blocked, 6 to influencers (>180,000 followers), 1 to an NGO with less than 2,000

1 followers and the last one to an individual (not influencer) with environmental concerns. As some images
2 could not be recovered as the accounts were blocked, the ranking was redone with those images that could be
3 recovered and eliminating those repeated after comparing their hash values. Most of the images (53)
4 belonged to various topics such as various objects, images from awareness campaigns, groups of people and
5 cartoons. Twenty of these images unequivocally captured trash in beach areas and the open sea. The rest of
6 27 images contained animals for raising awareness purposes and in fact, some of them were damaged or
7 entangled; the most common animal was the turtle (10), followed by fish (5), marine mammals (5), octopus
8 (4), birds (2) and crustaceans (1).

9 Another objective of this study was to automatically classify tags and images using available artificial
10 intelligence tools and explore their utility to inspect in a second phase the type and quantity of accumulated
11 marine litter in coastal areas. Unfortunately, the use of the Computer Vision API (v3.1; Azure Cognitive
12 Services) was deemed unsatisfactory in this study. From a random selection of 100 images, only 19 were
13 properly described by the tool. The tags attached by the software have not been helpful either. For example,
14 the image of a bucket full of cigarette butts on a sandy floor was described as “a bowl of nuts” and the tags
15 of this image were: ‘bowl’, ‘ground’, ‘floor’, ‘plate’ and ‘tea’. As stated in the API documentation, objects
16 are generally not detected if they are small (less than 5% of the image) and they are not detected if they are
17 arranged closely together, as in the case of hot spots of marine litter. This prospective study prevents us from
18 developing greater efforts, such as the training of a model for object detection, at least until there is a major
19 advance in this field of technology.

20 **4. Conclusions**

21 Twitter promotes the theory of public engagement, allowing users to have conversations, form communities,
22 share content and build relationships (Kietzmann et al., 2011). This paper advises of the potential of this
23 platform to create and spread environmental awareness, in this case, to combat the problem of marine litter in
24 the world connecting leaders, actors, companies, students and the public. This is the first time—at least to
25 the knowledge of the authors—that a scientific article explores the social network Twitter to analyze public
26 awareness about this issue.

27
28 The study describes a snapshot of an extremely dynamic social network spanning a period with an
29 exceptional pandemic situation world-wide. Most of the previous studies of various kinds that analyze
30 Twitter do so by only focusing on hashtags and keywords in English. The largest volume of tweets in our
31 dataset is in English, however and thanks to our multilingual approach, it is possible to analyze the
32 differences between countries from a broader point of view. The results show that countries such as the USA,
33 UK, Japan and Canada with a high population and digitization are those with the highest volume of tweets,

1 but when weighted by the number of active users then the topic is led by the UK and European countries like
2 Spain, France or Germany. In contrast, other high populated and digitized countries have relatively low
3 interest in this specific topic. The results also show the high engagement that occurs during the celebration of
4 “World Days” (e.g., #EarthDay, #WorldSeaTurtleDay or #PlasticFreeJuly) and related to the dissemination of
5 reports and scientific studies in traditional media such as newspapers or television which are echoed in this
6 SM. If we put the focus on a higher temporal resolution, we have found that the activity is mainly
7 concentrated on weekdays with differences between countries according to habits and work breaks.

8 Tweets are often informal, unstructured, making it challenging for deciphering a general discourse when read
9 individually. To know what the part of society that uses this social network is talking about, it is necessary to
10 analyze the tweets in an aggregate way. In this regard, the LDA technique has been effective in
11 distinguishing between different subtopics: i) “*Impact on wildlife*”, ii) “*Microplastics/Water pollution*”, iii)
12 “*Estimate of quantities/Reports*”, iv) “*Legislation/Protection*” and v) “*Recycling/Cleaning initiatives*”.
13 Besides, this technique has been useful to distinguish topics that were not directly associated with our
14 objective, allowing us to improve the cleaning technique of the original dataset. The impact of the COVID-
15 19 pandemic has also been evident in the messages with many of these messages referencing to mask and
16 glove waste and its impact on the environment. Within the topic “*Impact on wildlife*” the high number of
17 comments referring to entangled turtles is noteworthy. These tweets presented a slightly negative sentiment
18 as opposed to positive sentiment on the most general topic.

19 Our results show that NGOs, international organizations and academic institutions do not lead the
20 conversation on marine litter issues, in spite of their high research and environmental awareness efforts on
21 this topic. Bearing in mind the characteristics of the described snapshot may help to adjust the
22 communication plan in Twitter of those institutions that wish to play a relevant role to fight against marine
23 pollution by plastics, environmental awareness and scientific dissemination. This is relevant to offer citizens
24 reliable and certified information, as well as to change habits and to reinforce sustainable behaviour aimed at
25 protecting our seas. By identifying where it is tweeted from and in what language, institutions can focus their
26 efforts in those areas by combining, for example, several institutional Twitter accounts with a regional
27 perspective. Major events may impact how users discuss socio-scientific issues in online media. Thus, to
28 increase recruitment is useful to follow environmental “World Days” to identify and approach people in the
29 marine plastic pollution field with a more general public discourse, as well as invite influencers to join the
30 cause with an objective and truthful discourse. By identifying events and their participants, institutions can
31 increase and diversify their network and reach, useful to identify the type of audience they usually
32 communicate to.

33 To know the terminology used, the different sub-topics, feelings and reactions are useful clues/guidelines to

1 design an efficient communication strategy. The communication should be bidirectional considering that
2 Twitter is conversation and users choose who to follow. In addition, knowing favourite hours and days for
3 publication is relevant taking into account the “short lifespan” of a tweet (Wilson, 2019). Positive messages
4 are expected to reinforce recruitment and promote activism, for this reason, sentiment analysis is an
5 interesting approach to analyze the before and after of campaigns launched by an institution and that could
6 be followed under the same hashtag.

7 The images associated with a tweet contribute to increasing its impact; therefore, it is relevant to understand
8 what attracts the attention of the users of the platform. At this point, the use of images containing animals for
9 raising awareness purposes is a popular resource. Our results have also shown an increase in activity after the
10 dissemination of news about relevant or impacting scientific advances. Generally, it is believed that the more
11 individuals use social media, even if just to communicate and connect, the more likely they are to encounter
12 news (Huber et al., 2019). For this reason, institutions should be a source of scientific news that helps to
13 spread a truthful and contrasted discourse.

14 Another initial objective of this study was to verify if images that corresponded to waste in coastal areas
15 could be automatically filtered and used to create a map of coastal pollution by marine litter. However, the
16 artificial intelligence tool tested could not create correct descriptions of these images, among other reasons,
17 because the objects present were too small and appeared distorted. However, this social network —as well as
18 other popular ones in the use of images such as Instagram— have the potential to support local or regional
19 programs for coastal monitoring of marine litter, through the use of a specific hashtag or user mention. That
20 is why future studies should find the optimal way to use this social network to photograph coastal areas with
21 waste.

22 Finally, to contribute to this analysis over time, an interactive web application has been made available at
23 <http://twilitter.herokuapp.com/>. The tool allows the user to follow the temporal evolution, examine areas with
24 the highest volume of tweets, analyze sentiments or check the highest frequency of hashtags, among others.

25
26
27
28
29
30
31

1 **References**

- 2 Abreo, N.A.S., Thompson, K.F., Arabejo, G.F.P., Superio, M.D.A., 2019. Social media as a novel source of
3 data on the impact of marine litter on megafauna: The Philippines as a case study. *Marine Pollution*
4 *Bulletin* 140, 51–59. <https://doi.org/10.1016/j.marpolbul.2019.01.030>
- 5 Banda, J.M., Tekumalla, R., Wang, G., Yu, J., Liu, T., Ding, Y., Chowell, G., 2020. A large-scale COVID-19
6 Twitter chatter dataset for open scientific research -- an international collaboration.
7 *arXiv:2004.03688*.
- 8 Bank, M.S., Ok, Y.S., Swarzenski, P.W., 2020. Microplastic's role in antibiotic resistance. *Science* 369,
9 1315–1315. <https://doi.org/10.1126/science.abd9937>
- 10 Barker, J.L.P., Macleod, C.J.A., 2019. Development of a national-scale real-time Twitter data mining pipeline
11 for social geodata on the potential impacts of flooding on communities. *Environmental Modelling &*
12 *Software* 115, 213–227. <https://doi.org/10.1016/j.envsoft.2018.11.013>
- 13 Becken, S., Stantic, B., Chen, J., Alaei, A.R., Connolly, R.M., 2017. Monitoring the environment and human
14 sentiment on the Great Barrier Reef: Assessing the potential of collective sensing. *Journal of*
15 *Environmental Management* 203, 87–97. <https://doi.org/10.1016/j.jenvman.2017.07.007>
- 16 Blei, D.M., Ng, A.Y., Jordan, M.I., 2003. Latent dirichlet allocation. *J. Mach. Learn. Res.* 3, 993–1022.
- 17 Borrelle, S.B., Ringma, J., Law, K.L., Monnahan, C.C., Lebreton, L., McGivern, A., Murphy, E., Jambeck,
18 J., Leonard, G.H., Hilleary, M.A., Eriksen, M., Possingham, H.P., De Frond, H., Gerber, L.R.,
19 Polidoro, B., Tahir, A., Bernard, M., Mallos, N., Barnes, M., Rochman, C.M., 2020. Predicted
20 growth in plastic waste exceeds efforts to mitigate plastic pollution. *Science* 369, 1515–1518.
21 <https://doi.org/10.1126/science.aba3656>
- 22 Budenz, A., Klassen, A., Purtle, J., Yom Tov, E., Yudell, M., Massey, P., 2020. Mental illness and bipolar
23 disorder on Twitter: implications for stigma and social support. *J Ment Health* 29, 191–199.
24 <https://doi.org/10.1080/09638237.2019.1677878>
- 25 Carrasco-Polaino, R., Villar-Cirujano, E., Martín-Cárdaba, M.-Á., 2019. Redes, tweets y engagement:
26 análisis de las bibliotecas universitarias españolas en Twitter. *Profesional de la Información* 28.
27 <https://doi.org/10.3145/epi.2019.jul.15>
- 28 Clauset, A., Newman, M.E.J., Moore, C., 2004. Finding community structure in very large networks. *Phys.*
29 *Rev. E* 70, 066111. <https://doi.org/10.1103/PhysRevE.70.066111>
- 30 Collins, K., Shiffman, D., Rock, J., 2016. How Are Scientists Using Social Media in the Workplace? *PLOS*
31 *ONE* 11, e0162680. <https://doi.org/10.1371/journal.pone.0162680>
- 32 Cordasco, G., Gargano, L., 2011. Community Detection via Semi-Synchronous Label Propagation
33 Algorithms. *arXiv:1103.4550 [physics]*. <https://doi.org/10.1504/.045103>
- 34 Dahal, B., Kumar, S.A.P., Li, Z., 2019. Topic modeling and sentiment analysis of global climate change
35 tweets. *Soc. Netw. Anal. Min.* 9, 24. <https://doi.org/10.1007/s13278-019-0568-8>

- 1 Dang, T., Nguyen, N.V.T., Pham, V., 2018. HealthTvizier: Exploring Health Awareness in Twitter Data
2 through Coordinated Multiple Views, in: 2018 IEEE International Conference on Big Data (Big
3 Data). Presented at the 2018 IEEE International Conference on Big Data (Big Data), pp. 3647–3655.
4 <https://doi.org/10.1109/BigData.2018.8622445>
- 5 Duckett, P.E., Repaci, V., Duckett, P.E., Repaci, V., 2015. Marine plastic pollution: using community science
6 to address a global problem. *Mar. Freshwater Res.* 66, 665–673. <https://doi.org/10.1071/MF14087>
- 7 Eberhard, D., Simons, G.F., Fenning, C., 2019. *Ethnologue: Languages of the World*. SIL International.
- 8 European Commission, 2018. Single-use plastics: New EU rules to reduce marine litter [WWW Document].
9 Single-use plastics: New EU rules to reduce marine litter. URL
10 https://ec.europa.eu/commission/presscorner/detail/en/MEMO_18_3909 (accessed 12.28.20).
- 11 Filgueiras, A.V., Preciado, I., Cartón, A., Gago, J., 2020. Microplastic ingestion by pelagic and benthic fish
12 and diet composition: A case study in the NW Iberian shelf. *Marine Pollution Bulletin* 160, 111623.
13 <https://doi.org/10.1016/j.marpolbul.2020.111623>
- 14 Forleo, M.B., Romagnoli, L., 2021. Marine plastic litter: public perceptions and opinions in Italy. *Marine
15 Pollution Bulletin* 165, 112160. <https://doi.org/10.1016/j.marpolbul.2021.112160>
- 16 Fruchterman, T.M.J., Reingold, E.M., 1991. Graph drawing by force-directed placement. *Softw: Pract.
17 Exper.* 21, 1129–1164. <https://doi.org/10.1002/spe.4380211102>
- 18 Gago, J., Portela, S., Filgueiras, A.V., Salinas, M.P., Macías, D., 2020. Ingestion of plastic debris (macro and
19 micro) by longnose lancetfish (*Alepisaurus ferox*) in the North Atlantic Ocean. *Regional Studies in
20 Marine Science* 33, 100977. <https://doi.org/10.1016/j.rsma.2019.100977>
- 21 Gall, S.C., Thompson, R.C., 2015. The impact of debris on marine life. *Marine Pollution Bulletin* 92, 170–
22 179. <https://doi.org/10.1016/j.marpolbul.2014.12.041>
- 23 GESAMP, 2020. Proceedings of the GESAMP International Workshop on assessing the risks associated with
24 plastics and microplastics in the marine environment, Journal Series GESAMP Reports and Studies.
25 GESAMP Office, London.
- 26 GESAMP, 2015. Sources, fate and effects of microplastics in the marine environment: a global assessment
27 (No. 90), Rep. Stud. GESAMP. IMO, London.
- 28 Geyer, R., Jambeck, J.R., Law, K.L., 2017. Production, use, and fate of all plastics ever made. *Science
29 Advances* 3, e1700782. <https://doi.org/10.1126/sciadv.1700782>
- 30 Ghermandi, A., Camacho-Valdez, V., Trejo-Espinosa, H., 2020. Social media-based analysis of cultural
31 ecosystem services and heritage tourism in a coastal region of Mexico. *Tourism Management* 77,
32 104002. <https://doi.org/10.1016/j.tourman.2019.104002>
- 33 Goel, R., Sharma, R., 2020. Understanding the MeToo Movement Through the Lens of the Twitter, in: Aref,
34 S., Bontcheva, K., Braghieri, M., Dignum, F., Giannotti, F., Grisolia, F., Pedreschi, D. (Eds.), *Social
35 Informatics, Lecture Notes in Computer Science*. Springer International Publishing, Cham, pp. 67–

- 1 80. https://doi.org/10.1007/978-3-030-60975-7_6
- 2 Hagberg, A., Swart, P., S Chult, D., 2008. Exploring network structure, dynamics, and function using
3 networkx (No. LA-UR-08-05495; LA-UR-08-5495). Los Alamos National Lab. (LANL), Los
4 Alamos, NM (United States).
- 5 Heidbreder, L.M., Bablok, I., Drews, S., Menzel, C., 2019. Tackling the plastic problem: A review on
6 perceptions, behaviors, and interventions. *Science of The Total Environment* 668, 1077–1093.
7 <https://doi.org/10.1016/j.scitotenv.2019.02.437>
- 8 Huber, B., Barnidge, M., Gil de Zúñiga, H., Liu, J., 2019. Fostering public trust in science: The role of social
9 media. *Public Underst Sci* 28, 759–777. <https://doi.org/10.1177/0963662519869097>
- 10 Jambeck, J.R., Geyer, R., Wilcox, C., Siegler, T.R., Perryman, M., Andrady, A., Narayan, R., Law, K.L.,
11 2015. Plastic waste inputs from land into the ocean. *Science* 347, 768–771.
12 <https://doi.org/10.1126/science.1260352>
- 13 Karami, A., Lundy, M., Webb, F., Dwivedi, Y.K., 2020. Twitter and Research: A Systematic Literature
14 Review Through Text Mining. *IEEE Access* 8, 67698–67717.
15 <https://doi.org/10.1109/ACCESS.2020.2983656>
- 16 Kemp, S., 2020. Digital 2020: Global Digital Overview. Essential insights into how people around the world
17 use the internet, mobile devices, social media and ecommerce. We Are Social and Hotsuite®.
- 18 Kietzmann, J.H., Hermkens, K., McCarthy, I.P., Silvestre, B.S., 2011. Social media? Get serious!
19 Understanding the functional building blocks of social media. *Business Horizons*, SPECIAL ISSUE:
20 SOCIAL MEDIA 54, 241–251. <https://doi.org/10.1016/j.bushor.2011.01.005>
- 21 Law, K.L., Starr, N., Siegler, T.R., Jambeck, J.R., Mallos, N.J., Leonard, G.H., 2020. The United States’
22 contribution of plastic waste to land and ocean. *Science Advances* 6, eabd0288.
23 <https://doi.org/10.1126/sciadv.abd0288>
- 24 Letierce, J., Passant, S., Brelisn, J.G., 2010. Understanding how Twitter is used to spread scientific message.
25 Presented at the Proceedings of the WebSci10: Extending the Frontiers of Society On-Line.
- 26 Li, X., Xie, Q., Huang, L., Yuan, Z., 2017. Twitter Data Mining for the Social Awareness of Emerging
27 Technologies, in: 2017 Portland International Conference on Management of Engineering and
28 Technology (PICMET). Presented at the 2017 Portland International Conference on Management of
29 Engineering and Technology (PICMET), pp. 1–10. <https://doi.org/10.23919/PICMET.2017.8125279>
- 30 Makita, M., Mas-Bleda, A., Morris, S., Thelwall, M., 2021. Mental Health Discourses on Twitter during
31 Mental Health Awareness Week. *Issues in Mental Health Nursing* 42, 437–450.
32 <https://doi.org/10.1080/01612840.2020.1814914>
- 33 Manikonda, L., Beigi, G., Liu, H., Kambhampati, S., 2018. Twitter for Sparking a Movement, Reddit for
34 Sharing the Moment: #metoo through the Lens of Social Media. arXiv:1803.08022 [cs].
- 35 Martínez-Rojas, M., Pardo-Ferreira, M. del C., Rubio-Romero, J.C., 2018. Twitter as a tool for the

1 management and analysis of emergency situations: A systematic literature review. *International*
2 *Journal of Information Management* 43, 196–208. <https://doi.org/10.1016/j.ijinfomgt.2018.07.008>

3 Media Insight Project, 2017. ‘Who shared it?’ How Americans decide what news to trust on social media.
4 American Press Institute. URL [https://www.americanpressinstitute.org/publications/reports/survey-](https://www.americanpressinstitute.org/publications/reports/survey-research/trust-social-media/)
5 [research/trust-social-media/](https://www.americanpressinstitute.org/publications/reports/survey-research/trust-social-media/) (accessed 5.3.21).

6 Mehrotra, R., Sanner, S., Buntine, W., Xie, L., 2013. Improving LDA topic models for microblogs via tweet
7 pooling and automatic labeling, in: *Proceedings of the 36th International ACM SIGIR Conference on*
8 *Research and Development in Information Retrieval, SIGIR '13*. Association for Computing
9 Machinery, New York, NY, USA, pp. 889–892. <https://doi.org/10.1145/2484028.2484166>

10 Milani, E., Weitkamp, E., Webb, P., 2020. The Visual Vaccine Debate on Twitter: A Social Network Analysis.
11 *Media and Communication* 8, 364–375. <https://doi.org/10.17645/mac.v8i2.2847>

12 Mitrano, D.M., Wohlleben, W., 2020. Microplastic regulation should be more precise to incentivize both
13 innovation and environmental safety. *Nature Communications* 11, 5324.
14 <https://doi.org/10.1038/s41467-020-19069-1>

15 Mohammadi, E., Thelwall, M., Kwasny, M., Holmes, K.L., 2018. Academic information on Twitter: A user
16 survey. *PLOS ONE* 13, e0197265. <https://doi.org/10.1371/journal.pone.0197265>

17 Napper, I.E., Thompson, R.C., 2020. Plastic Debris in the Marine Environment: History and Future
18 Challenges. *Global Challenges* 4, 1900081. <https://doi.org/10.1002/gch2.201900081>

19 Newman, M.E.J., 2010. *Networks: an introduction*. Oxford University Press, Oxford ; New York.

20 Pabortsava, K., Lampitt, R.S., 2020. High concentrations of plastic hidden beneath the surface of the Atlantic
21 Ocean. *Nature Communications* 11, 4073. <https://doi.org/10.1038/s41467-020-17932-9>

22 Parton, K., Galloway, T., Godley, B., 2019. Global review of shark and ray entanglement in anthropogenic
23 marine debris. *Endang. Species. Res.* 39, 173–190. <https://doi.org/10.3354/esr00964>

24 Plackett, R., Kaushal, A., Kassianos, A.P., Cross, A., Lewins, D., Sheringham, J., Waller, J., von Wagner, C.,
25 2020. Use of Social Media to Promote Cancer Screening and Early Diagnosis: Scoping Review. *J*
26 *Med Internet Res* 22, e21582. <https://doi.org/10.2196/21582>

27 Polinario, J., 2016. *Cómo divulgar ciencia a través de las redes sociales*, Investigación. Circulo Rojo.

28 Retka, J., Jepson, P., Ladle, R.J., Malhado, A.C.M., Vieira, F.A.S., Normande, I.C., Souza, C.N., Bragagnolo,
29 C., Correia, R.A., 2019. Assessing cultural ecosystem services of a large marine protected area
30 through social media photographs. *Ocean & Coastal Management* 176, 40–48.
31 <https://doi.org/10.1016/j.ocecoaman.2019.04.018>

32 Röder, M., Both, A., Hinneburg, A., 2015. Exploring the Space of Topic Coherence Measures, in:
33 *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM*
34 *'15*. Association for Computing Machinery, New York, NY, USA, pp. 399–408.
35 <https://doi.org/10.1145/2684822.2685324>

- 1 Ruiz-Frau, A., Ospina-Alvarez, A., Villasante, S., Pita, P., Maya-Jariego, I., de Juan, S., 2020. Using graph
2 theory and social media data to assess cultural ecosystem services in coastal areas: Method
3 development and application. *Ecosystem Services* 45, 101176.
4 <https://doi.org/10.1016/j.ecoser.2020.101176>
- 5 Sayyadiharikandeh, M., Varol, O., Yang, K.-C., Flammini, A., Menczer, F., 2020. Detection of Novel Social
6 Bots by Ensembles of Specialized Classifiers. *Proceedings of the 29th ACM International*
7 *Conference on Information & Knowledge Management* 2725–2732.
8 <https://doi.org/10.1145/3340531.3412698>
- 9 Shimkhada, R., Attai, D., Scheitler, A.J., Babey, S., Glenn, B., Ponce, N., 2021. Using a Twitter Chat to
10 Rapidly Identify Barriers and Policy Solutions for Metastatic Breast Cancer Care: Qualitative Study.
11 *JMIR Public Health and Surveillance* 7, e23178. <https://doi.org/10.2196/23178>
- 12 Sievert, C., Shirley, K., 2014. LDAvis: A method for visualizing and interpreting topics, in: *Proceedings of*
13 *the Workshop on Interactive Language Learning, Visualization, and Interfaces*. Association for
14 *Computational Linguistics*, Baltimore, Maryland, USA, pp. 63–70. [https://doi.org/10.3115/v1/W14-](https://doi.org/10.3115/v1/W14-3110)
15 [3110](https://doi.org/10.3115/v1/W14-3110)
- 16 Smith, A., 2015. “Wow, I didn’t know that before; thank you”: How scientists use Twitter for public
17 engagement. *Journal of Promotional Communications* 3.
- 18 STATISTA, 2020. Leading countries based on number of Twitter users as of October 2020. URL
19 <https://www.statista.com/statistics/242606/number-of-active-twitter-users-in-selected-countries/>
20 (accessed 12.28.20).
- 21 Steinskog, A., Therkelsen, J., Gambäck, B., 2017. Twitter Topic Modeling by Tweet Aggregation, in:
22 *Proceedings of the 21st Nordic Conference on Computational Linguistics*. Association for
23 *Computational Linguistics*, Gothenburg, Sweden, pp. 77–86.
- 24 Taddy, M.A., 2011. On Estimation and Selection for Topic Models. arXiv:1109.4518 [stat].
- 25 Teoh, D., Shaikh, R., Vogel, R.I., Zoellner, T., Carson, L., Kulasingam, S., Lou, E., 2018. A Cross-Sectional
26 Review of Cervical Cancer Messages on Twitter during Cervical Cancer Awareness Month. *J Low*
27 *Genit Tract Dis* 22, 8–12. <https://doi.org/10.1097/LGT.0000000000000363>
- 28 The Pew Charitable Trusts and SYSTEMIQ, 2020. *Breaking the Plastic Wave: A Comprehensive Assessment*
29 *of Pathways Towards Stopping Ocean Plastic Pollution*. The Pew Charitable Trusts and SYSTEMIQ.
- 30 Van Noorden, R., 2014. Online collaboration: Scientists and the social network. *Nature News* 512, 126.
31 <https://doi.org/10.1038/512126a>
- 32 Vince, J., Hardesty, B.D., 2017. Plastic pollution challenges in marine and coastal environments: from local
33 to global governance. *Restoration Ecology* 25, 123–128. <https://doi.org/10.1111/rec.12388>
- 34 Vraga, E.K., Stefanidis, A., Lamprianidis, G., Croitoru, A., Crooks, A.T., Delamater, P.L., PFOSER, D.,
35 Radzikowski, J.R., Jacobsen, K.H., 2018. Cancer and Social Media: A Comparison of Traffic about

1 Breast Cancer, Prostate Cancer, and Other Reproductive Cancers on Twitter and Instagram. *Journal*
2 *of Health Communication* 23, 181–189. <https://doi.org/10.1080/10810730.2017.1421730>

3 W3Techs, 2020. Usage statistics of content languages for websites. URL
4 https://w3techs.com/technologies/overview/content_language (accessed 12.28.20).

5 Wadhwa, V., Latimer, E., Chatterjee, K., McCarty, J., Fitzgerald, R.T., 2017. Maximizing the Tweet
6 Engagement Rate in Academia: Analysis of the AJNR Twitter Feed. *American Journal of*
7 *Neuroradiology* 38, 1866–1868. <https://doi.org/10.3174/ajnr.A5283>

8 Warner, K., Linske, E., Mustain, P., Valliant, M., Leavitt, C., 2020. Choked, Strangled, Drowned: The
9 Plastics Crisis Unfolding In Our Oceans. OCEANA.

10 Wilson, C., 2019. Updated: Lifespan of a Social Media Post. Max Influence. URL
11 <https://mtomconsulting.com/updated-lifespan-of-a-social-media-post/> (accessed 12.28.20).

12 Xue, J., Chen, J., Hu, R., Chen, C., Zheng, C., Su, Y., Zhu, T., 2020. Twitter Discussions and Emotions About
13 the COVID-19 Pandemic: Machine Learning Approach. *Journal of Medical Internet Research* 22,
14 e20550. <https://doi.org/10.2196/20550>

15 Yoosefi Nejad, M., Delghandi, M.S., Bali, A.O., Hosseinzadeh, M., 2019. Using Twitter to raise the profile
16 of childhood cancer awareness month. *Netw Model Anal Health Inform Bioinforma* 9, 3.
17 <https://doi.org/10.1007/s13721-019-0206-4>

18 Zhou, S., Kan, P., Huang, Q., Silbernagel, J., 2021. A guided latent Dirichlet allocation approach to
19 investigate real-time latent topics of Twitter data during Hurricane Laura. *Journal of Information*
20 *Science* 01655515211007724. <https://doi.org/10.1177/01655515211007724>

21 Zubiaga, A., Ji, H., 2014. Tweet, but verify: epistemic study of information verification on Twitter. *Soc.*
22 *Netw. Anal. Min.* 4, 163. <https://doi.org/10.1007/s13278-014-0163-y>

23

24

25

26

27

28

29

30

31

32

33

34

35

1 **List of tables**

2 Table 1. Example of some fields from tweets in the dataset. Time and original message were directly
 3 obtained from Twitter API. Clean text after processing the original message is also shown in the table. The
 4 city, state, and country fields were calculated from user profile and added to the dataset. Polarity [-1, 1] and
 5 subjectivity [0, 1] are also shown. User identifier is not shown to protect the identities of the Twitter users.

6

Time (UTC)	Original message	Clean text	City, State, Country	(Polarity, subjectivity)
Sat Aug 08 19:52:52 2020	A pocos metros de una playa del sur.. recogimos en 10 minutos 7 bolsas de latas y botellas de plástico abandonadas...falta concienciación de limpieza y civismo en la isla. @GranCanariaCab https://t.co/A4pCuWyP2G	meters southern beach collected minutes bags abandoned cans plastic bottles lack awareness cleanliness civility island	Las Palmas de Gran Canaria, Islas Canarias, Spain	(-0.1, 0.05)
Tue Jul 28 22:34:55 2020	Last week, we did a beach cleanup with @LagunaOceanFdn at Aliso Beach, the end of the Aliso Creek watershed! We picked up 10 pounds of trash in ONE hour, prevented gulls from eating plastic cups, saw native &... https://t.co/hUpG2Vu03B	last week beach cleanup aliso beach end aliso creek watershed picked pounds trash one hour prevented gulls eating plastic cups saw native	Costa Mesa, California, USA	(0, 0.07)
Thu Aug 13 10:45:04 2020	Nearly half of the plastic found in the ocean comes from fishing nets. People are reducing the consumption of plastics, but given that the scientific community warned that by 2050 there will be more plastic in the ocean than fish, it is not enough. https://t.co/CmDTOyhBoL	nearly half plastic found ocean comes fishing nets people reducing consumption plastics given scientific community warned plastic ocean fish enough	Allerdale, England, UK	(0.11, 0.38)

7
8
9
10
11
12
13
14
15
16
17

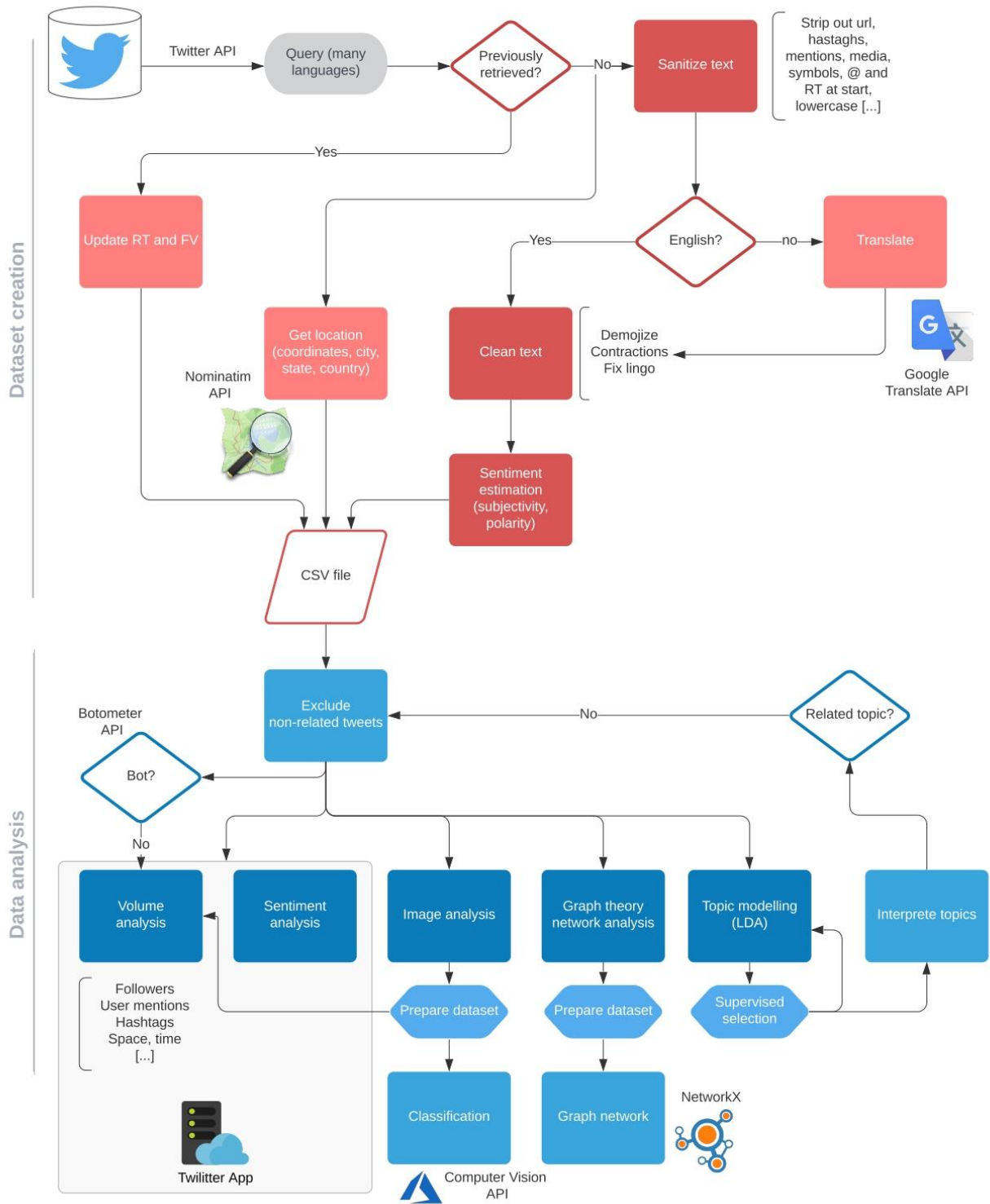
1 Table 2. The 10-top most salient words for 6 topics generated from LDA analysis. Words are in descending
 2 order of relevance following the definition of Sievert et al. (2014), computed with a weight parameter $\lambda = 1$
 3 and $\lambda = 0.4$.
 4

Human interpretation	Top-10 salient terms ($\lambda = 1$)	Top-10 salient terms ($\lambda = 0.4$)	Selection of terms ($\lambda < 0.4$)
Impact on wild life	bag, waste, turtle, use, life, straws, stop, animals, trash, single	turtle, bag, animals, straws, waste, stop, use, life, killing, fishing	birds, creatures, jellyfish, dolphins, killing, nets, harming, mammals, entangled, whales, stomach
Microplastics / Water pollution	water, bottle, new, pollution, clean, bag, micro, study, food, litter	water, bottle, micro, atlantic, new, particles, tiny, times, study, surface	butts, cigarettes, cans, technology, discovered, clothes, fragments, fibers
Estimate of quantities / Reports	pollution, million, year, waste, tons, world, use, fish, people, end	million, year, tons, pollution, estimated, metric, likely, lets, report, waste	lockdown, coronavirus, consequence, biodiversity, mediterranean
Legislation / Protection	global, pollution, help, hi, states, climate, protecting, members, legislation, treaty	global, hi, states, protecting, climate, member, legislation, treaty, requires, encouraging	unenvironment, truth, warming, loveplanet, protectdepends, spain, fuels, fosil, oil, forest, melting
Recycling / Cleaning initiatives	free, pollution, help, single, use, people, clean, save, july, make	free, july, solution, act, minute, challenge, cleaner, helps, communities, signed	Trump, refuse, congress, fund, movement
Gorillaz*	gorillaz, album, days, like, masks, good, best, demon, song, love	gorillaz, album, days, masks, demon, best, song, face, covid, good	-

5 *Gorillaz topic (British virtual band), completely unrelated to marine pollution by plastics and microplastics. LDA revealed this topic
 6 and helped to improve cleaning in the dataset.
 7
 8
 9
 10
 11
 12
 13
 14
 15
 16

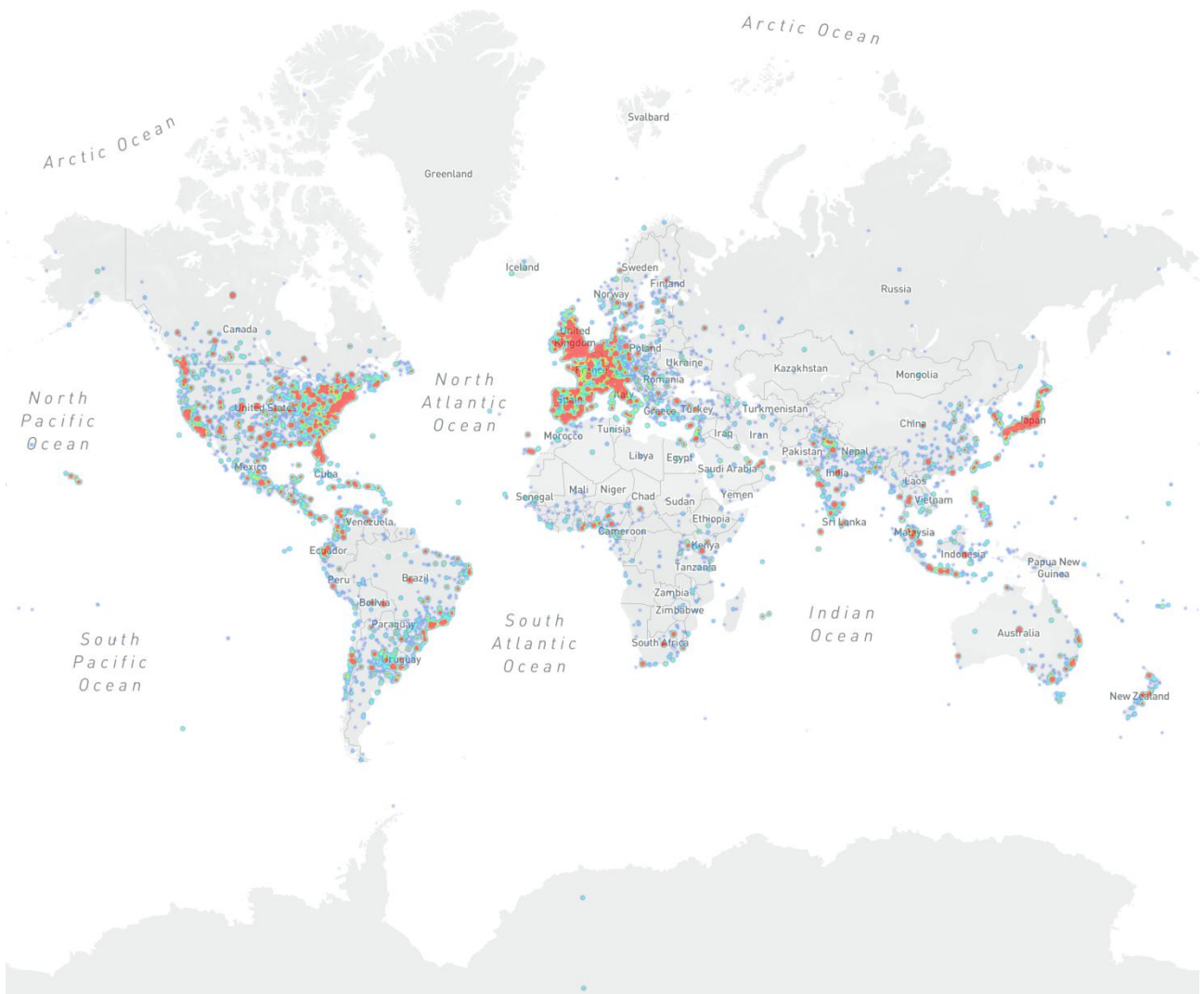
1 List of figures

2 Figure 1: Flowchart of the dataset preparation and data analysis.



3

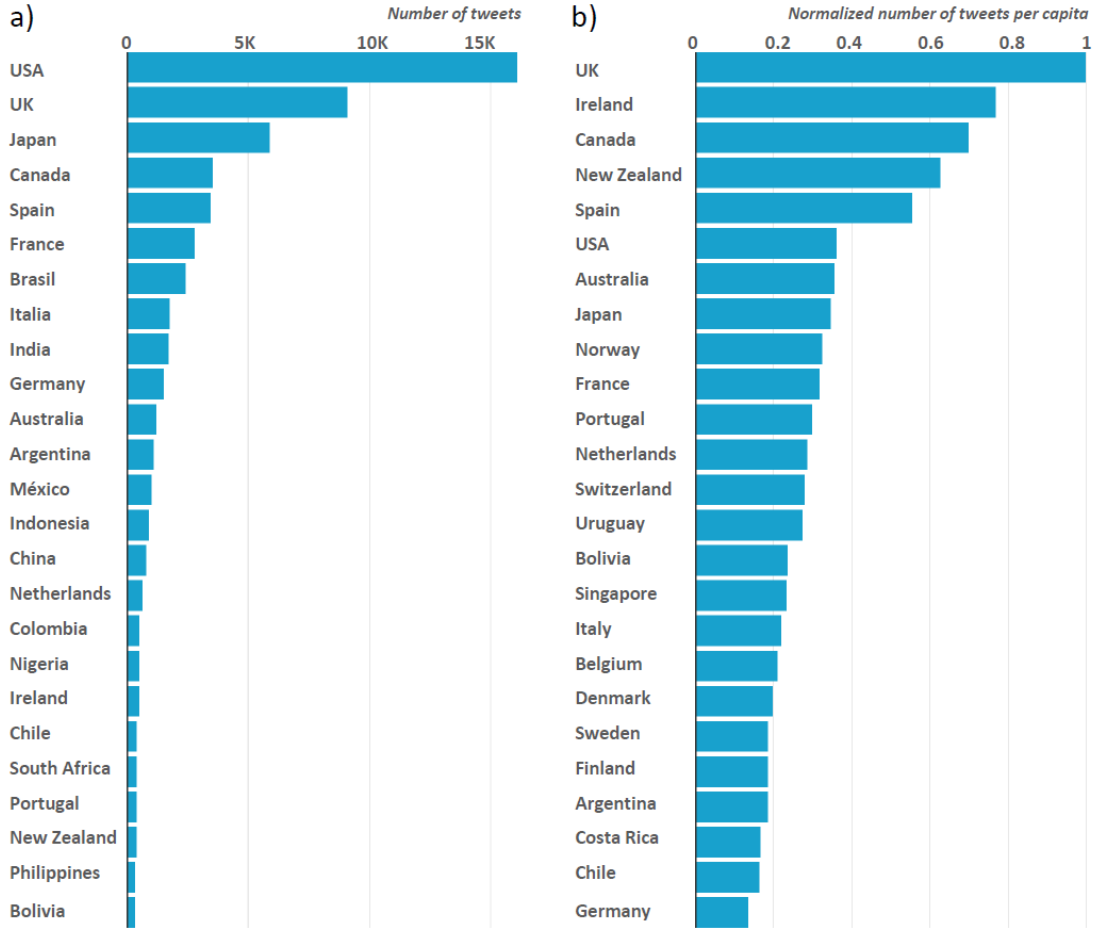
1 Figure 2: Heatmap of tweets from 19th March to 26th November 2020. Locations were retrieved from both
2 geotagged tweets and from Twitter user's profile.
3



4
5
6
7
8
9
10
11
12
13
14
15

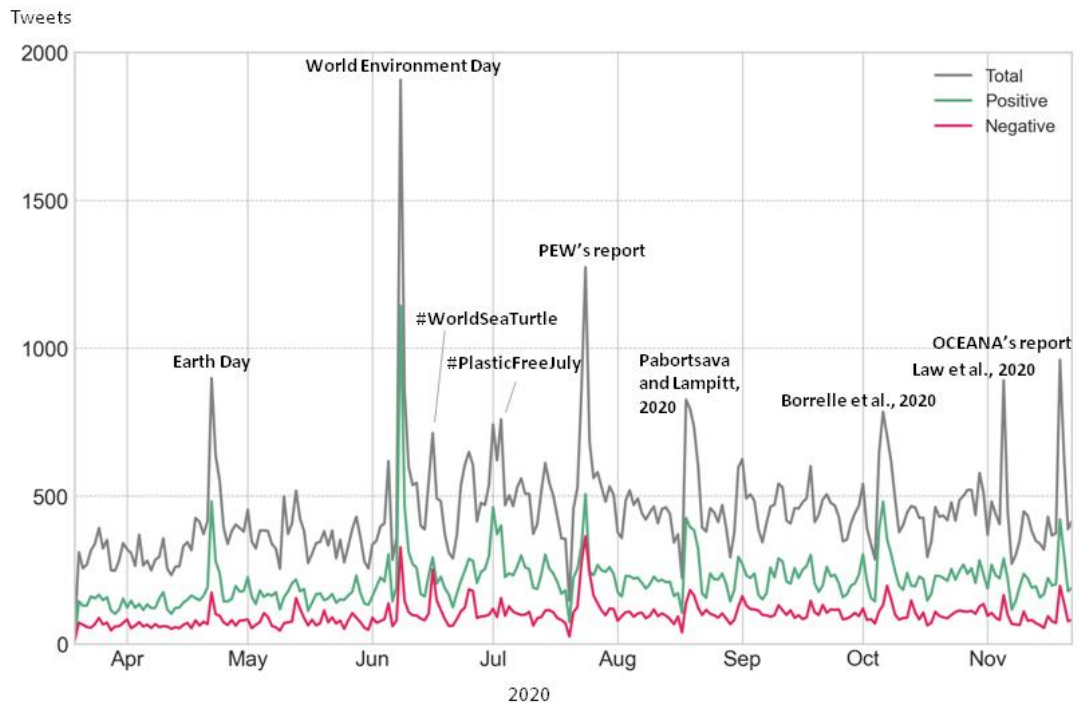
1 Figure 3: Top-25 countries sorted by a) number of tweets and b) number of tweets per capita normalized by
 2 the maximum ratio among countries (UK).

3

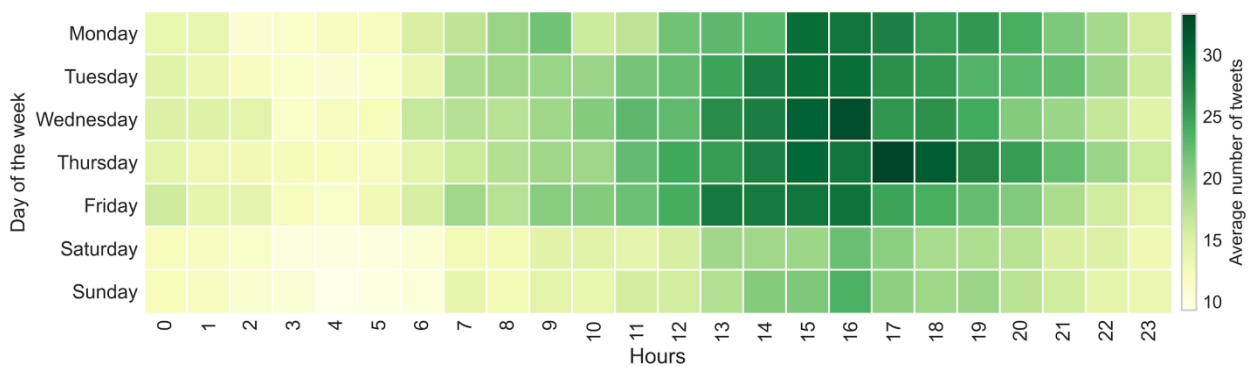


4
 5
 6
 7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17

- 1 Figure 4: Total number of tweets per day (dark grey), with positive sentiment (green) and negative (red).
- 2 Relevant events are annotated on the figure.

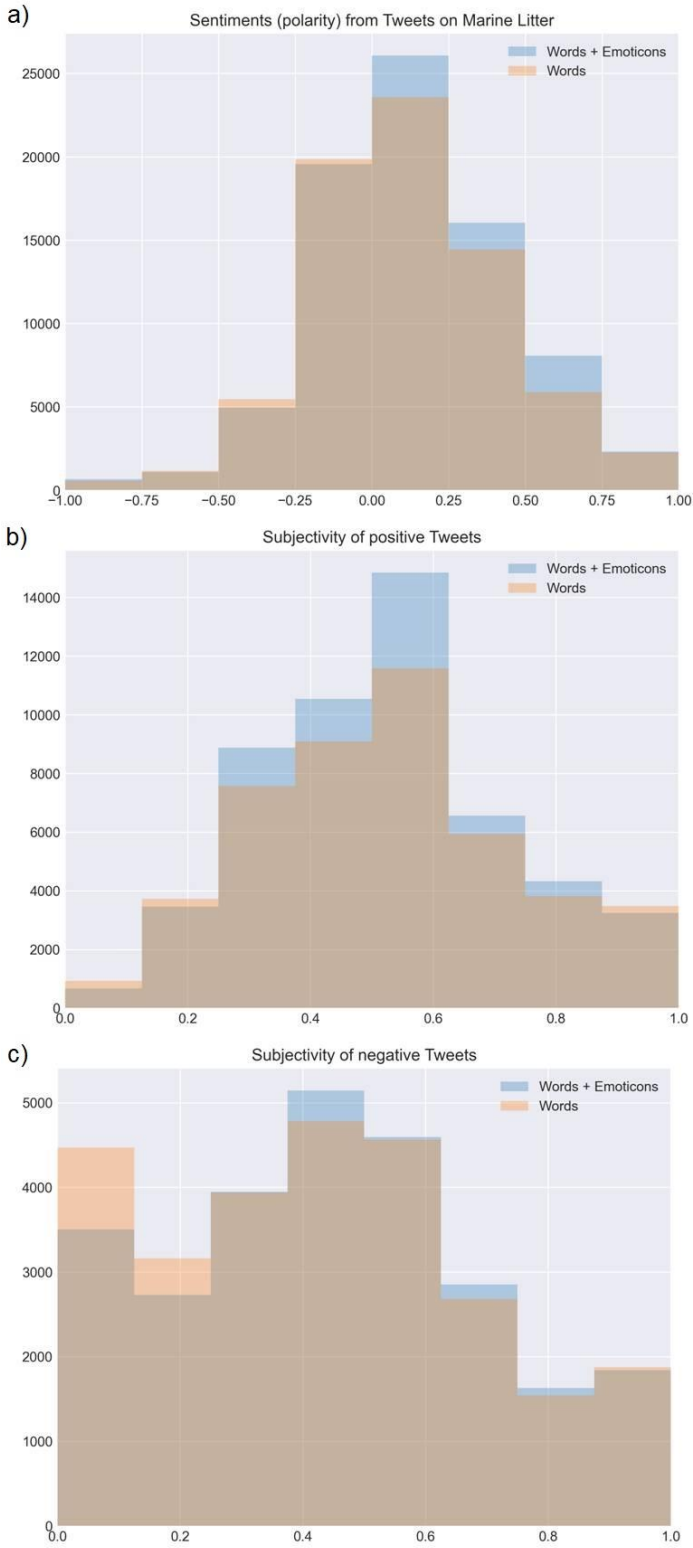


- 3
- 4 Figure 5: Average weekly tweets within each time slot (UTC hours)



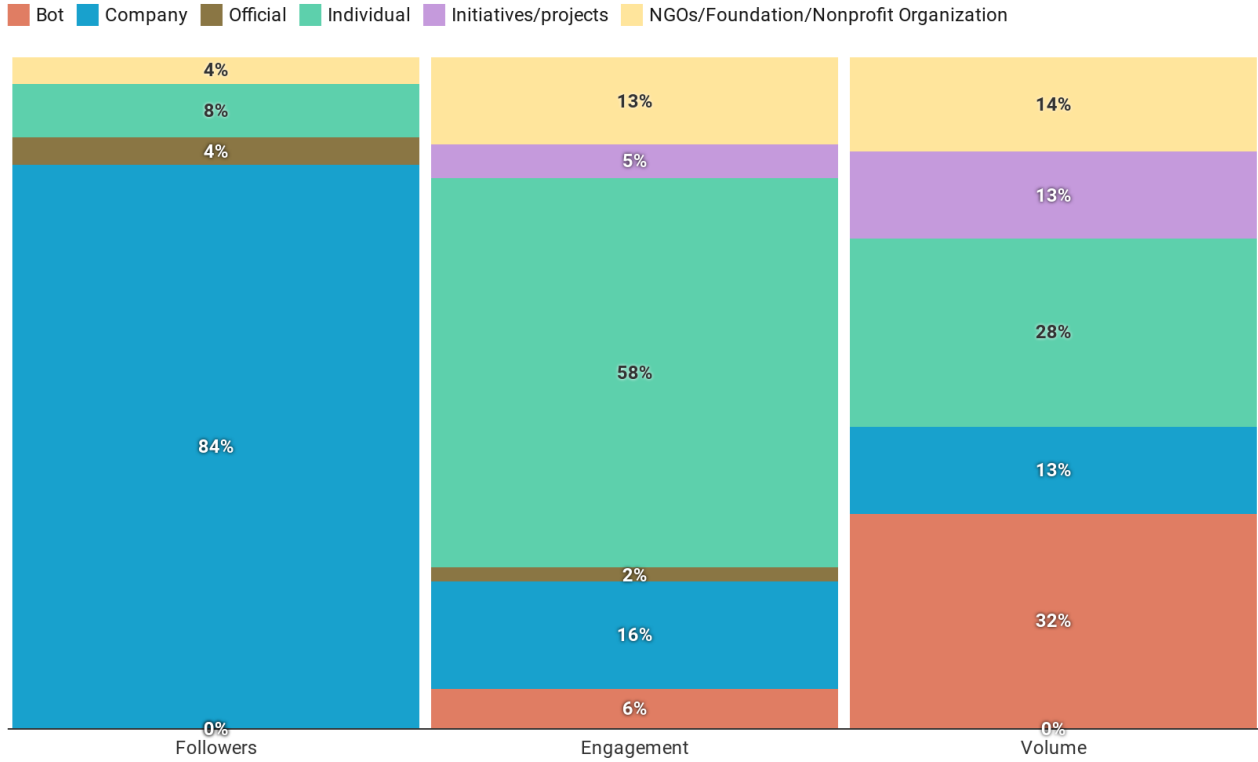
- 5
- 6
- 7
- 8
- 9
- 10

1 Figure 6: Sentiment analysis distribution for the total of tweets: a) polarity of positive and negative tweets, b)
2 subjectivity of positive tweets and c) subjectivity of negative tweets. Subjectivity ranges from 0 (very
3 objective) to 1 (very subjective).



4
5

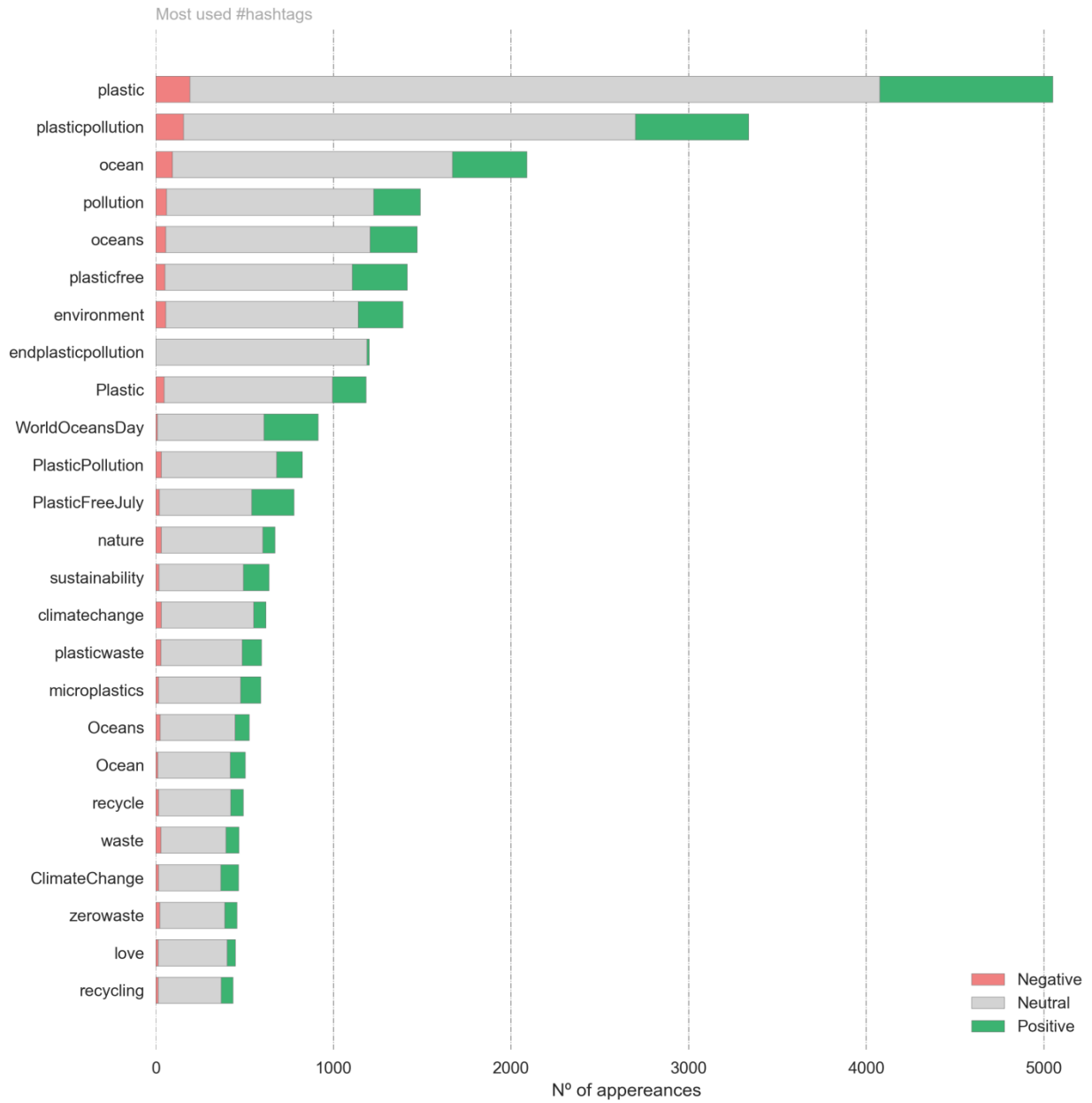
1 Figure 7: Top-100 users by tweet volume, engagement (likes plus retweets) and the number of followers
 2 categorized by bots, companies, official institutions, initiatives/projects and NGOs/foundation/Nonprofit
 3 organizations. In the analysis by followers, the category of companies is made up of 90% by mass media.



4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

1 Figure 8: Top-25 most used hashtags sorted by the number of appearances in the dataset. The proportion of
 2 negative (<-0.3), positive (>0.3) and neutral sentiment is also shown.

3



4

5

6

7

8

9

10

11

