

Towards the operational use of tweets data in high impact weather scenarios: data mining and analytics in Basque Country

S. Gaztelumendi 1,2

1- Basque Meteorology Agency (EUSKALMET), Parque tecnológico de Álava, Miñano, Araba, Basque Country. 2- TECNALIA BRTA (Basque Research and Technology Alliance), Meteorology Area, Parque tecnológico de Álava, Miñano, Araba, Basque Country.

Abstract

In this contribution, we present different aspects related to the operational use of Twitter data in the context of high impact weather scenarios at local level. We present some results and experiences from a proof of concept project that demonstrate how different data mining and advanced analytics techniques, can be used in order to include social media information for different operational tasks and particularly during severe weather events in a weather service context like the Basque Meteorology Agency.

Introduction

Although social media industry is nowadays a very congested Marketplace, Twitter continues to maintain its status as a popular social media platform with 330 million monthly active users and 145 million daily active users sending more than 6,000 tweets every second in the world (Businessofapps 2021). In Spain case, 85% population are social media users, with around 5 Spanish profiles for a population around 47 million (Statista 2021). In the autonomous country of Basque Country (CAE), with 2.17 million inhabitants in 7,234 km², we can estimate that around 70% of total population are social media users with a Twitter penetration of around 20%, after some inference from Spanish data on a study of Basque Youth and Social Media (EUSTAT 2020, Basque Youth Observatory 2018). It is noteworthy that Twitter is a social tool that enables users to post messages (tweets) of up to 280 characters supporting a wide variety of social communication practices including photo and video attachment. The Basque Meteorology Agency (Euskalmet) with more than 115.3 K followers remain as one of the most popular accounts in Basque Country (Gaztelumendi et al. 2013).

Methodology and Data

In this paper we summarize our experience during a proof of concept project for automatic real time Twitter mining and analysis and the key aspects to consider for the development of an operational tool for Twitter API data exploitation at local level. We present the main challenges and problems that we have had to deal with, including how to deal with the lack of geolocation information, since in the case of the Basque country, as in other parts of the world, tweets containing geotags are the exception, not the rule.

Key findings and conclusions

- The increase in social interest around weather during severe weather episodes, is clearly reflected in social networks, despite the size of the territory (approx. 100kx100km) and the inherent limitations of the Twitter API (theoretically only 1% of the live stream can be collected for free) (see fig 4)
Twitter content examination and analysis in order to extract some derived quantities aggregated for different spatial and temporal scales could be the basis for an automatic surveillance and monitoring social system that could be useful on real or deferred time.
It is important to note that different Knowledge aggregation levels are present from spatial, temporal and content characteristics and that information extraction is also possible at different levels of aggregation (see fig 4 and 5)
Not only text content could be useful, available metadata (user characteristics, attached information, etc) could be an important resource.
Exact location information for a particular tweet (geo coordinates from where a tweet is send) is rarely available but in many cases must be inferred from location field at different level of spatial aggregations (municipalities, county, historical territories, etc.) (see fig2 and fig 3)
As a general rule, each single tweet without any location information or outside the area of surveillance or content of interest (Basque country and surroundings) is not considered (see fig 4)
All tweets are categorized for potential credibility based on users characteristics. Known users are categorized according to their general reliability considering different characteristics (official sources, media, weather regular contributors, etc.). Unknown users are considered at low level credibility until an analysis is made. As a general rule each single tweet content is considered inaccurate and suspicious.
A plausible rule-based methodology could be implemented, for tweets text content and location mining, at a relatively low cost and serves as the basis for further and more complex developments.
In spite of the general positioning of Euskalmet in the Basque Country and the high penetration of @Euskalmet in the Basque Twitter community, the officially used severe weather hashtags (defined in Basque language and theoretically the key for discussion and topic analysis) apparently have a limited impact outside the institutional users.
Dictionaries / lexicons (topics, locations, etc.) must be implemented taking into account particular local idiomatic aspects (Basque and Spanish mixed words) and peculiarities of language usage in social networks (abbreviations, spelling mistakes, etc.)
Different metrics (e.g. increasing rates of number of a particular topic tweets) could be real time monitored at different spatial and temporal scales as a sort of "social sensor" network. Population distribution and other socio-cultural aspects need to be included in order to extract conclusions in the side of impact.
In this PoC, project, a preliminary real-time tweet collection and classification system is implemented, as well as some reduced monitoring and surveillance tools.
Visual data analytics techniques are essential for rapid human interpretation and they might be actionable at real time.
Both, fully automatic and supervised systems are needed for a full operational exploitation of available data.

Acknowledgements

The authors would like to thank the Department of Security of the Basque Government and particularly to the Directorate of Emergencies and Meteorology for operational service financial support. We also would like to thank all our colleagues from DAEM, EUSKALMET, AZTI and TECNALIA for their daily effort in promoting valuable research and services for the Basque Society. We would also like to thank Twitter users, R community and all institutions and people that support open data and tools. This work has been partially funded by the LIFE-URBANKLIMA2050 project



Contact info: santiago.gaztelumendi@tecnalia.com
TECNALIA
Parque Tecnológico de Bizkaia C/ Geldo Edificio 700 E-48160 DERIO (Bizkaia) Spain www.tecnalia.com

Results and Discussion

General Objective: To have better general impact information before/during after severe events that makes it possible to deal correctly with different Euskalmet's real time (RT) and delayed time (DT) operations in severe weather scenarios or other natural hazards events.
Final Goal: The implementation of an operational surveillance system of social media, news and emergencies data, fully integrated with actual monitoring and open source capabilities present in Twitter.
Context: Small country, highly concentrated population, bilingual with different penetration level of Basque Language. High usage of internet, connected mobile phones, and social network usage with an acceptable proportion of Twitter usage. Free Twitter API availability before/during after severe events that makes it possible to deal correctly with different Euskalmet's real time (RT) and delayed time (DT) operations in severe weather scenarios or other natural hazards events.
Final Goal: The implementation of an operational surveillance system of social media, news and emergencies data, fully integrated with actual monitoring and open source capabilities present in Twitter.
Context: Small country, highly concentrated population, bilingual with different penetration level of Basque Language. High usage of internet, connected mobile phones, and social network usage with an acceptable proportion of Twitter usage. Free Twitter API availability before/during after severe events that makes it possible to deal correctly with different Euskalmet's real time (RT) and delayed time (DT) operations in severe weather scenarios or other natural hazards events.

Twitter, that originally was designed for effective and efficient two way communication, could be considered today as one of the simplest and most redundant public communications tool. Twitter provides high-volume, high-velocity and high-variety unstructured data (big data) that can be used to support decision-making (e.g. O'Leary, 2015), particularly in the Euskalmet case (Gaztelumendi et al 2015b) and around meteorological business (Gaztelumendi et al 2016c).

There is substantial quantitative and qualitative information available for mining and analysis in Twitter, including number of tweets, number of retweets, number of followers and many other statistics and metrics (e.g. O'Leary 2015). In addition, there is substantial non-numeric qualitative information in terms of text that we need to automatically convert in some kind of actionable quantitative data.
Twitter mining (e.g. O'Leary 2015) and Twitter analytics (e.g. Kumar et al 2013) is concerned with providing structure to the unstructured data in order to extract and support information. Twitter text messages and metadata are our "data mine" and we mine that data for its potential usage in the field of local impact weather. For this purposes different modules and submodules are implemented for data acquisition, filtering, cleaning, geolocation extraction, topic classification, analysis of content and knowledge exploitation.
When dealing with content analysis, different approaches are possible (e.g. Kruspe et al 2021, Reuter et al 2017, Kuman et al 2013). In our case, different and specific tools were prepared for different topics and in two different flows: one for Twitter API querying and other for general message content mining. In this PoC, just fully human-based lexicons are prepared with the most frequent terms related with the so defined key subjects nearly the same as causes available in the warning/alert/alarm operational system (OV 2018, Gaztelumendi et al 2012). Lexicons are defined as bilingual, containing Basque as Spanish terms. Note that single queries to the Twitter API are limited to 500 characters so, if needed, multiples queries are performed.

Different typologies could be used to categorize social media messages, based on content and metadata, covering many dimensions as information provided, emotional content, source, credibility, time, location, etc. (see table 4.1 in chapter 4 Castillo 2016 for more details). In our case, we include all metadata directly available for each tweet from the API (more than 90 fields), including time, geolocation (if present), and user id. We also include some relevant derived information from the different classification modules. In this PoC, mainly based on location and severe weather typology according to warning/alert/alarm operational criteria.
Such modules are executed sequentially in an automatic process that organizes and understands the large collections of available twitter data, by assigning "tags" or categories according to each individual tweet characteristics (mainly from field "location" and "text"). As a first iteration we implement different rule-based classification modules by directly programming a set of hand-made rules based on the content of textual fields. Defined rules are able to extract location information from the "location" and "text" fields and to discern between tweets of different topics by looking directly at semantically relevant elements of a particular "text" field content. Rules are defined using different lexicons (lists of words) and using different metrics for similarities as the highest frequency specific topics words.

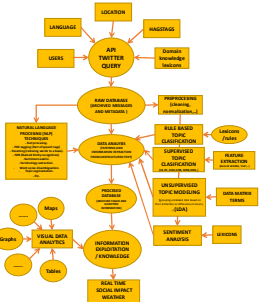


Fig. 1. General concept framework.

Data provision module. Set of Twitter APIs queries for data acquisition based on different scripts and tweet package executed recursively each 5 minutes. Queries are pre-filtered (API) by key-words and location. Final datatables (API) are conveniently merged and labeled.

User classification module. Extract relevant user information for categorization purposes. Users are classified and labeled according to different categories. In the PoC just 4 users typologies are included: governmental/public, news/media, collaborators, others. In next steps further classification must be included with the final idea of characterizing potential credibility to a given content using derived information available from the user side (e.g. influence metrics). User related labels are included for each tweet in the tweets DF.

Location classification module. Location data extraction at different aggregation level, directly from geotag/available data (exact punctual lon lat coordinates, box lon, etc) or indirectly extracting information from the text location field (municipality, county, province, etc.). Those tags oriented to identifying the presence of key terms from different geographical lexicons prepared for this purpose. All the information available is included in the tweets DF.

Message content classification module. Categorization of tweets messages with rule-based methodology (e.g. Beria 2020). Different keyword-based rules are used (in a simple system) to separate messages into categories defined by the presence of different key-words in the message. Although it is generally considered an ineffective approach, is the simplest and most straightforward methodology, and this may work for certain information categories that have a small, well-defined, unambiguous set of terms that are highly discriminative (Castillo, 2016), as is the case. This module includes various Natural Language Processing techniques for text pre-processing (e.g. tokenization, stop words extraction, etc.).

Representation module. At the end of the automatic process, messages are converted to a format suitable for automatic representation. Final DF contains original data and new features from text and metadata content analysis. In addition, visual data analytics techniques are used to achieve a better understanding of the data. In this PoC, just some basic graphics are generated.

Fig. 2. Schematic representation of PoC modules.

Twitter and its structure.

A "tweet" is a general message posted on Twitter.com. It is restricted to 280 characters. Though most tweets contain media, it is possible to embed links (URLs), pictures, GIFs, GIFs, videos or emojis. Once a tweet is sent by a user it becomes immediately visible to all followers. If it is found useful it can be retweeted to other users. A retweet ("RT" for short) is a tweet by a user that has been retweeted by user 1 to all of user 1's followers. Thus, retweeting is a way of measuring how popular the tweet is.

This occurs either as a "re-tweet", where the entire tweet is re-forwarded following by one of the secondary followers, or as a "retweeted" tweet, where a part of the tweet is forwarded or what is known as "quote retweet" (where the original tweet content is "quoted" by the "RT" symbol). A tweet is shared through the original author's handles or "handles" (the "RT" symbol is transmitted). The original Twitter handles may be mentioned in any subsequent tweets, so it is usually feasible to track the original message. Users or handles are recognized by the "RT" symbol. A user can direct a message to another user by adding the handle, with the "@" symbol. An @, or a mention, is when you include somebody's identifier in the tweet. The person will be alerted that you mentioned them. It is used to send a public, "hey, over here," or to add someone on a conversation that isn't currently happening.

Tweets usually contain components called "hashtags" which are words that capture the subject of the tweet. They are prefixed by the "#" character. Hashtags are clickable, too, so you can tap on a hashtag to see all the tweets related to that topic.

A user can "favorite" a tweet (paraphrase to a "like" on Facebook or Instagram). A reply to a tweet means responding to a message or tweet from a person visible to retweets to broadcast (like forwarding an email) tweet message by one reply to another.

There are two ways to reply to tweets. There is a "reply" where you "G@vernate" in a tweet. You can use a DM (direct message) which is sent only to the recipient (a private email).

Twitter APIs. Twitter APIs provide access to a variety of different resources, including information about tweets, user profiles, trends, and more. There are two main types of APIs: the REST API and the Streaming API. The REST API is used for fetching tweets and user information, while the Streaming API is used for real-time data.

Twitter users. Twitter users are individuals or organizations that have created a Twitter account. They can post tweets, follow other users, and interact with tweets. Twitter users are categorized into different groups based on their activity and influence.

Twitter content. Twitter content refers to the text, images, and videos that are posted on the platform. This content is used to communicate information and engage with other users. Twitter content is analyzed using various techniques to extract insights and trends.

Twitter location. Twitter location data is used to track the geographic distribution of tweets. This data is used to identify trends and patterns in user behavior across different regions and countries.

Twitter analytics. Twitter analytics provide insights into the performance of tweets and user engagement. This data is used to optimize marketing campaigns and improve user experience.

Twitter trends. Twitter trends are popular topics or keywords that are currently being discussed on the platform. These trends are used to identify current events and public opinion.

Twitter hashtags. Twitter hashtags are used to categorize tweets and make them easier to find. They are used to track specific topics and trends on the platform.

Twitter retweets. Twitter retweets are a way for users to share and amplify content. They are used to increase the visibility of tweets and engage with a larger audience.

Twitter mentions. Twitter mentions are used to tag and mention other users in tweets. They are used to acknowledge and interact with other users on the platform.

Twitter direct messages. Twitter direct messages are private messages sent between users. They are used for one-on-one communication and are not visible to other users.

Twitter search. Twitter search is used to find tweets and users based on specific keywords and filters. It is used to track trends and identify relevant content on the platform.

Twitter API. The Twitter API is a set of tools that allow developers to interact with the Twitter platform. It is used to build applications and services that integrate with Twitter.

Twitter data. Twitter data is the collection of all tweets and user information on the platform. It is used for research and analysis to understand user behavior and trends.

Twitter trends. Twitter trends are popular topics or keywords that are currently being discussed on the platform. These trends are used to identify current events and public opinion.

Twitter analytics. Twitter analytics provide insights into the performance of tweets and user engagement. This data is used to optimize marketing campaigns and improve user experience.

CONTEXT

Social and cultural aspects in Basque Autonomous Community (CAE). The Basque language is one of the two official languages in the territory of CAE spoken by more than 2 million of inhabitants, along with Spanish language spoken by nearly all population. Basque is a minority general language that has received different information and standardization historical processes (e.g. Urrutia et al 2020), nowadays it benefits from international language status.

Basque Language. The Basque language is one of the two official languages in the territory of CAE spoken by more than 2 million of inhabitants, along with Spanish language spoken by nearly all population. Basque is a minority general language that has received different information and standardization historical processes (e.g. Urrutia et al 2020), nowadays it benefits from international language status.

Population Spatial Distribution. The total population of the Basque Country Autonomous Community (CAE) (see Fig. 4) is around 2,170,000 with a mean density of 384 inhabitants/km². 80% of the population of CAE resides in urban areas, particularly in municipalities with more than 10,000 inhabitants (see Fig. 4). The distribution pattern is unequal (see Fig. 4) with 52% in the territory of Bizkaia, 25% in Gipuzkoa and 23% in Araba. In Araba, 76.2% of the population resides in one of the three capital cities. The population of Araba is highly concentrated, with almost 75% living in Vitoria-Gasteiz (see Fig. 4). In Bizkaia, the population is almost 70% live around great Bilbao area (see Fig. 4). In Gipuzkoa, the population is more spread out than in the other territories, with nearly 50% living in small or medium municipalities with less than 10,000 inhabitants (see Fig. 4). See EMS2020-2021-0011 https://www.tecnalia.com

Basque Language. The Basque language is one of the two official languages in the territory of CAE spoken by more than 2 million of inhabitants, along with Spanish language spoken by nearly all population. Basque is a minority general language that has received different information and standardization historical processes (e.g. Urrutia et al 2020), nowadays it benefits from international language status.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

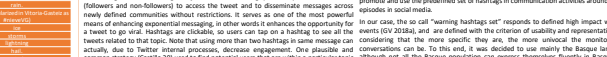
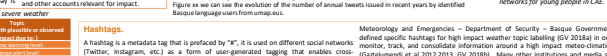
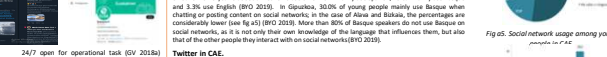
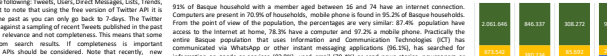
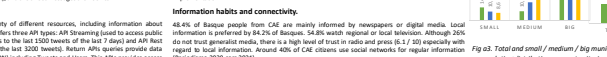
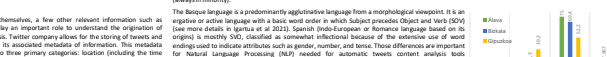
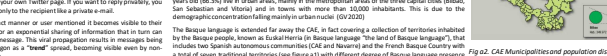
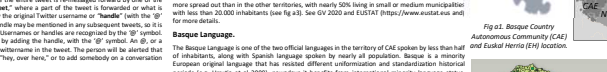
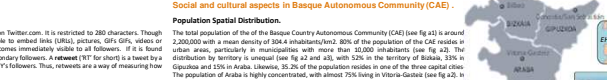
Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.

Information fields and connectivity. 48% of Basque people from CAE are mainly informed through digital media, local television and radio. 48% of Basque people from CAE are mainly informed through digital media, local television and radio.



References

List of references including works by Aguirre, Beria, Castillo, and others, covering topics like social media analysis, weather forecasting, and language use.