

# Event detection on a Twitter dataset

**Tanguy Albrici**

tanguy.albrici@epfl.ch

**Davide Di Dio**

davide.didio@epfl.ch

**Bruno Wicht**

bruno.wicht@epfl.ch

## Abstract

Social networks are widely used to discuss real-life regional and international events. In this report, we describe a hashtag-based way to perform event detection and localization from a dataset of tweets geolocated in Switzerland. The method proposed is effective for major incidents and popular festivities, and is suitable for events lasting one or multiple days. Our procedure detected a total of 7082 distinct events and found 440 meaningful local events locations.

## 1 Dataset Description and Pre-processing

At first, we extract the hashtags from each tweets. Each tweet is stored in a dataframe that only contains the informations we require, namely:

- The tweet id
- The user id
- The longitude and latitude
- The hashtags extracted from the text
- The day, month and year of the tweet creation

After the hashtag extraction, we only keep the tweets with at least one hashtag. This dataframe is called `df_tag` in our notebook.

## 2 Data Manipulation

### 2.1 Grouping by hashtag

In order to implement event detection efficiently, we need to compute a dictionary containing for each hashtags the ids of each post that contains it. This is done in the function `group_by_hashtag(...)`. This function iterates over every post in the database and adds the

id of the post in the dictionary entry of the hashtag. In this method, we also compute the number of unique authors for this hashtag. With this value, we can easily filter out hashtags that were tweeted by only a few users.

## 3 Data analysis and visualization

This section covers how we visualize our data and allow us to get a feel on how to detect events.

### 3.1 Visualizing hashtag frequency

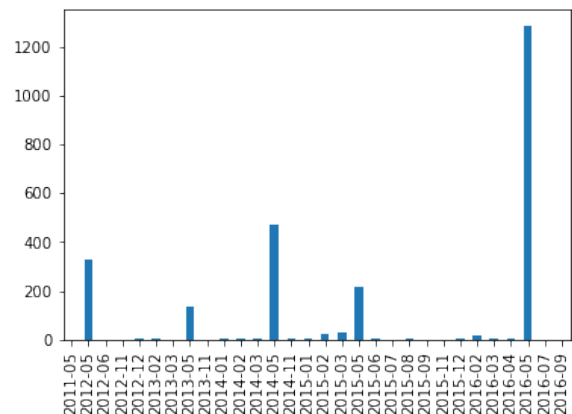


Figure 1: For the top months, plot the number of tweets with #eurovision

Before starting the event detection part, it is important to get a feel of how the number of tweets typically varies for some given hashtag. In Figure 1, we see a plot of the top 30 months that had the highest number of tweets with the hashtag #eurovision. Note that this plot doesn't necessarily contain consecutive months, but only the more important ones, in chronological order. We can very clearly see an increase in the number of tweets during the month of May of each year. Indeed, the Eurovision contest takes place during that month, and it is obvious that people talk more about it while it is happening. We will be using

this knowledge to develop an algorithm to detect those sudden spikes.

### 3.2 Visualizing event localisation

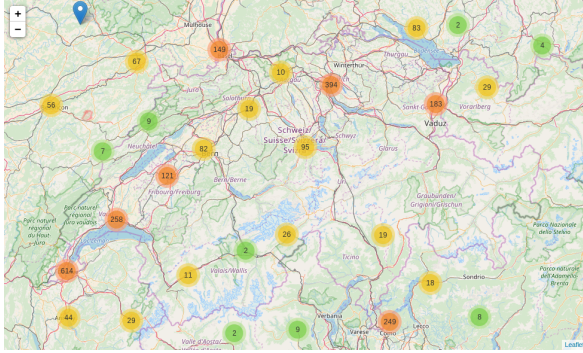


Figure 2: Map with every tweet for #eurovision

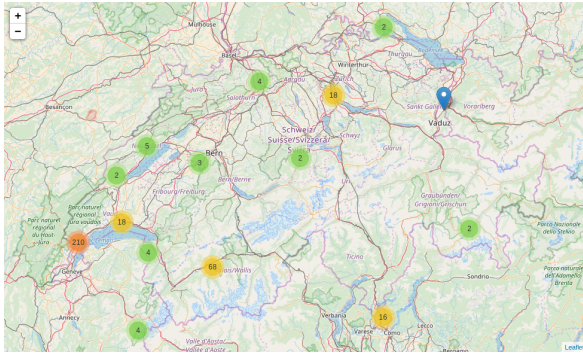


Figure 3: Map with every tweet for #paleo

The map are created using folium and are displayed in order to get an intuition on the geographic repartitions of tweets. We use a Marker-Cluster to group every tweet close in space.

In the map shown in Figure 2, we see that there are more tweets in the cities but overall the distribution is quite uniform if we consider the population density. On the other hand, in the map shown in Figure 3, we see that most of the tweets are located in Nyon, which is indeed where the Paleo music festival takes place. Hence, we can see that local events lead to a concentration of tweets around the position of the event, whereas international events have a much more uniform geographic tweet distribution. We will use this knowledge to find the type of an event, and estimate its location in case of a local event.

## 4 Event detection

In this part of the project we focus on automatic event detection.

### 4.1 Filtering out irrelevant hashtags

In the data manipulation, we constructed a dictionary containing each hashtag and their respective number of post and unique authors. Since hashtags that do not have enough posts and that are posted by only a few users are not likely to be considered events by our detection algorithm. We will be filtering out every hashtag that have less than 100 overall posts and less than 50 total authors.

### 4.2 Main Challenges

There were two main challenges we faced when trying to come up with an accurate event detection algorithm.

The first one is that our dataset spans almost 7 years (from early 2010 to late 2016) and thus the frequency of tweets is obviously much higher at the end of our dataset than at the beginning, as shown in Figure 4. Therefore, we had to compare the popularity of a hashtag at a certain date to some kind of baseline, so that our results can adapt to the increasing number of tweets.

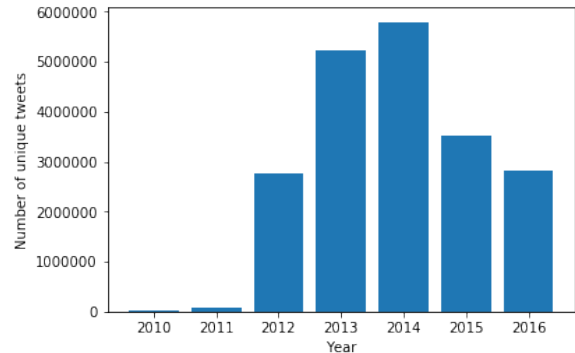


Figure 4: Number of tweets per years in our dataset. Note that the years 2010 and 2016 are not included in full in the dataset.

The second challenge is that there are multiple types of events, and they can have very different characteristics. For example, for an event that is scheduled to happen at a certain date, the number of tweets about it gradually increases up to that date. Whereas for a sudden unpredictable event, nobody talks about it and then a big peak happens at the time of the event. Moreover, some events might last multiple days or even weeks, as opposed to one day. Therefore this makes it harder to find an algorithm that detects well all types of events.

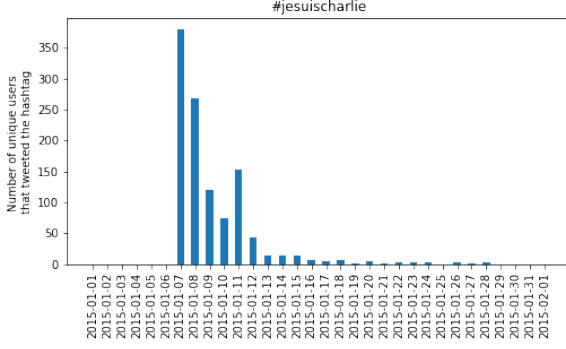


Figure 5: Example of an unpredictable event

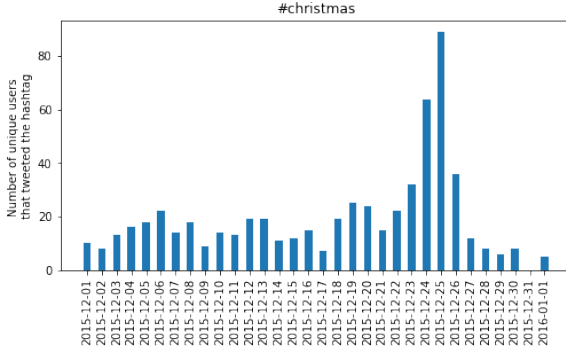


Figure 6: Example of a predictable event

### 4.3 Event Detection Algorithm

#### 4.3.1 Algorithm setup

First of all, we decided to use as a metric of hashtag popularity the number of unique users that tweeted that hashtag. We think this is more relevant than using the total number of tweets, because it prevents us from detecting events based on tweets that were posted by a few users only (e.g. twitter bots or extremely active users). For simplicity, we also decided to detect only the day when an event happens, and not the time, as it would be more complicated. Thus, our algorithm needs as input only the number of unique authors that tweeted a given hashtag, for each day spanned by our dataset, and for each hashtag.

#### 4.3.2 Algorithm description

The main idea of our event detection algorithm is to compute what we will call an *event score*, for each day spanned by our dataset and for each hashtag. Then, we apply a threshold on that event score to get dates that are detected as an event. The event score is obtained by dividing some absolute measure of popularity, by a baseline. That baseline is obtained by computing an average of the number of unique users that tweeted the hash-

tag over some odd number of days  $N$  around the day for which we compute the event score. Dividing by this baseline solves the first problem we talked about. The absolute measure of popularity is obtained by a local weighted average over some odd number of days  $M$  around the day for which we compute the event score, where  $M \ll N$ . More specifically, if the number of unique users that tweeted a hashtag at a certain day  $d$  is given by  $U_d$ , and that the kernel for the local weighted average is given by

$$K = [w_{\lfloor -\frac{M}{2} \rfloor} \dots w_0 \dots w_{\lfloor \frac{M}{2} \rfloor}]$$

then, the event score at day  $d$  is obtained by :

$$event\_score(d) = \frac{\sum_{i=\lfloor -\frac{M}{2} \rfloor}^{\lfloor \frac{M}{2} \rfloor} w_i * U_{d+i}}{\frac{1}{N} \sum_{j=\lfloor -\frac{N}{2} \rfloor}^{\lfloor \frac{N}{2} \rfloor} U_{d+j}}$$

The kernel for the weighted local average should be some kind of gaussian, so that the maximum weight is for day  $d$ . The aim of this local weighted average is to better detect anticipated events, or events lasting a few days. This was implemented to address the second challenge mentioned earlier.

#### 4.3.3 Grouping events in time

After running the algorithm described above, we have a list of hashtags, and the corresponding days where the event score was above the threshold. What we do next is group this list of days for each hashtag so that the dates that are very close are represented as one event over multiple days. More specifically, we grouped two dates together if they were at most 2 days apart.

#### 4.3.4 Parameters used

To detect events on the Swiss tweets dataset, we used the following parameters:

$$N = 35$$

$$M = 3$$

$$K = [0.1 \ 0.8 \ 0.1]$$

$$threshold = 4$$

## 5 Event Localization

Now that we have a list of events, we want to determine the location of an event based on the latitude/longitude data from the tweets that caused this event to be detected. To do so, we chose to compute the median of the latitude and longitude of the tweets.

But for event that are not necessarily physically happening at some place or that are just trends, a location might not be relevant. Thus in order to determine if the event location is any good, we also compute the mean standard deviation between the computed event location and the location of the tweets. Then we decide that an event is local if the mean standard deviation is below a threshold. If it is not, the event location computed is most likely not relevant and thus we drop it.

## 6 Results

Here are some statistics to quantify our results:

- Out of 6197 different hashtag on which we run the event detection algorithm, we found events for 2108 of them.
- By grouping by similar dates, we then found 7028 different events.
- Out of these 7028 events, we found a relevant location for 440 of them.

### 6.1 Event Detection

We correctly detected the events we took as example previously, mainly #jesuischarlie and #christmas, as we can see in Figure 7 and 8. Actually, for #christmas, we detected events on the 24th and 25th of every year between 2012-2015, which is quite good (Note that december 2016 was not included in our dataset).

### 6.2 Event Localization

A map containing all 440 events for which a location was found is shown in Figure 9. Unsurprisingly, we can see that most of the events happen in cities, they are especially numerous in Geneva and Zürich.

## 7 Conclusion

To conclude, we have implemented an event detection algorithm, that has shown to detect many events and produce good results. Moreover, we have managed to estimate the location of the events for which the tweets were very localized.

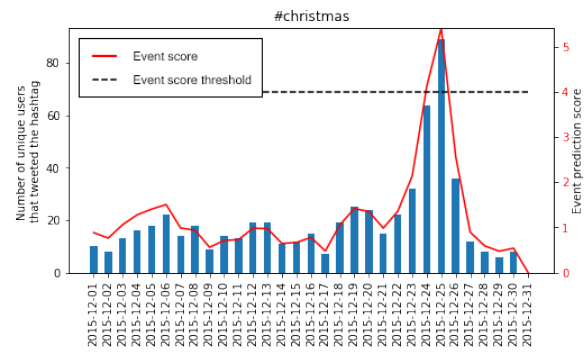


Figure 7: Computed event score (in red) for the event #christmas

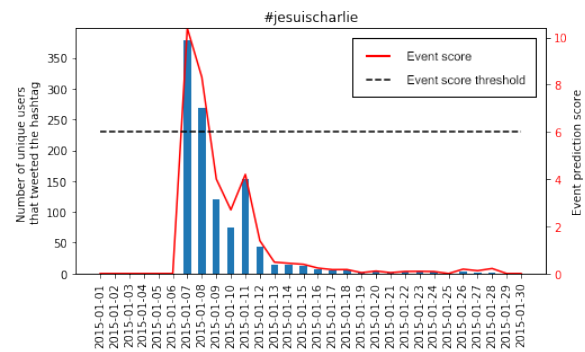


Figure 8: Computed event score (in red) for the event #jesuischarlie

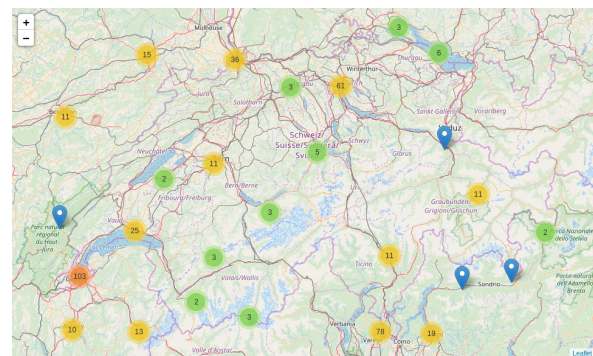


Figure 9: Map of all localized events