

Spleetter: a PACT-based Twitter Statistical Analyzer

Matteo Lissandrini
University of Trento
ml@disi.unitn.eu

Davide Mottin
University of Trento
mottin@disi.unitn.eu

ABSTRACT

Twitter is the biggest micro-blogging system and users therein produce every day a 175 million of short messages¹, namely tweets. The ability to perform analyses fast on the tweets is not only useful, but also vital to retrieve news and real-time statistics. However, since the number of different statistics to compute on the same data is big (e.g., the hashtag trend over the time, the number of tweets per user), there is the need of performing several operations at the same time. On one hand we can consider a Map-Reduce system in order to split the operations into different machines running several map and reduce steps. On the other hand, a map and reduce paradigm may not suffice. We propose to solve this problem with a PACT flow, that is a semantically richer map-reduce like-system. Our solution is both simple and highly modular, being composed by simple operations that can be easily parallelised. We show how to compute a set of interesting statistics out of 22 million tweets downloaded in a period of two weeks, showing the polarity of each tweet and filtering the data to produce interesting analyses.

Keywords

Twitter, Stratosphere, Data Analysis, Sentiment Analysis

1. INTRODUCTION

The ability to perform various statistical analyses with Twitter has attracted much interest and research in the last few years [2, 3] and used to influence politics [9] and advertising [1]. With a 140,000,000 subscribers to its service, Twitter is far and away the leading social network for the real time sharing and re-sharing of information, becoming the most used communication media in the spreading of socio-cultural events across the world²[4]. Moreover, Twitter has an enormous advantage over other social networks like Facebook in one key area: while people on Facebook tend to

¹http://www.mediabistro.com/alltwitter/twitter-stats_b32050

²<http://www.umpf.co.uk/blog/?p=6830>

friend their friends, people on Twitter tend to follow their interests³. Previous works as [4], [6] and [5] present different techniques to identify peak of activities around particular keywords, *hashtags* or users. They highlight possibilities for topic discovery techniques for informative summarisation of events, and for user-friendly visualisations of specific streams of postings. All these previous approaches aim to do real time identification of data peaks, smart labelling or even post-mortem analysis of events. We would like to mine a real time stream of posts and to identify trends and patterns that can allow us to forecast which topics, keywords or events will become popular in the near future. To do this we need to actually mine social media timelines to collect relevant statistics and insights on how events resonates in a social network, but the huge amount of information per day cannot be processed by a single machine. .

To this end, we propose a schema based on the PACT paradigm implemented in Stratosphere, that is an enhanced map-reduce system able to mix second order functions. Our goal is to conduct an analysis of a big dataset, automatically downloaded using the Twitter streaming APIs. We propose a PACT flow or program⁴ that takes in input a set of tuples, having tweets and user ids, a set of users and a vocabulary and produces several statistics to be of use to the data analyst in order to perform in depth data analysis and event discovery. We now describe the basic structure of the Twitter system and an highlight of the proposed solution discussed in Section 2.

1.1 Twitter structure

Twitter is a micro-blogging system designed to allow users to send short messages having a maximum of 140 characters, called *tweets*. In the text, users are allowed to specify *hashtags* that are sequence of characters usually describing an argument and marked by the character '#'. Users can also reference other users with '@user' notation. A particular kind of reference is a *retweet* which is a tweet preceded by the tag "RT" and the user name that first posted the message. A user is *retweeting* the tweet from another user when she thinks that it is worth spreading such post to a broader audience. The last information in the tweets are the urls, that are usually shortened using some available URL shortener.

1.2 Proposed solution

³<http://dcurt.is/twitters-graph>

⁴here flow or program are used interchangeably

We aim to find information regarding topic trends, analysing user post trends, polarity of the tweets with respect to the timing of those tweets, and with respect also to hashtags contained in them. The set of operations we propose to implement using the PACT programming is the following:

- **Tweet Cleansing:** we take in input the tuples containing the tweets and a dictionary of english words and we filter out hashtags, user mentions and tweets having a number of english words less than a threshold. The cleaned data are then use throughout the rest of the flow.
- **Polarity extraction:** we use a well-known library for sentiment analysis to extract the general polarity of each tweet. The polarity is defined as value between $[-5, 5]$, where a positive number means that the user is talking about something in a favorable manner. Conversely, if a text has a polarity close to -5 the words in the text are dissenting. From [8] we will adopt the SentiStrength classifier, which was built especially to cope with sentiment detection in short informal text. It combines a lexicon-based approach with more sophisticated linguistic rules.
- **Hashtag analysis:** we produce a various statistics about hashtags. Those take into account the time in which the hashtag appeared and disappeared (i.e. is not used anymore), the maximum and minimum number of mentions per hashtag with the timestamp of those peaks,
- **Topic Analysis:** we perform an analysis on positive and negative trends per topic, aggregating for each hashtag, for each point in time, the average polarity expressed in tweets containing it. We compute also the average Emotional Divergence[7] for each tag.

The document is structured as follows. Section 3 describes the dataset we downloaded and used in our experiments. We propose a solution and describe the PACT program in detail in Section 2. We also show the results of the analysis, with increasing size of the dataset in Section 4. Concluding, in Section 5 we describe the problems and the issues we encountered during the development of the solution and we remark our findings.

2. SOLUTION

We describe here the PACT program diagram we implemented to produce the statistics on tweets that we want. To recap briefly, a PACT program is a generalisation of the Map-Reduce paradigm, in which sequence of second order functions are issued in parallel or in sequence and combined to execute complex tasks. The programming paradigm defines 5 second order functions: map, reduce, match, cross and co-group. It allows the user to specify any kind of combination between them, in any order. We propose a flow that uses all the operator except the cross. For ease of explanation we describe the PACT program as composed in two different blocks: data cleaning (Section 2.1) and the computation of the statistics (Section 2.2).

2.1 Data cleaning

In the data cleaning part we take as input a comma separated file having the format $\langle tweet_id, user_id, tweet \rangle$. Table 1 shows an excerpt of the tweet tuples.

Note that it is not always easy to find interesting information from arbitrarily short-texts. Therefore, we propose a first flow that tries to remove useless or uninformative tweets by filtering-out the ones that do not contain a minimum portion of english words. Figure 1 depicts the sequence of operations we designed to clean the data in the preliminary phase.

Once we loaded the tweets into tuples we clean them, from hashtags, user-mentions and URLs. Tweets are then used in two separate flows: (1) we split them into words in order to count the english words and (2) we perform a sentiment analysis over the text, in order to evaluate the positive and negative polarity expressed by them, as described in Section 1. Note that the SentiStrength library is well suited to analyse informal text, so we will apply this evaluation on the whole text, before stemming and further cleaning.

In order to restrict the search space to those tweets we consider relevant, we import a dictionary of english words and we count english words in each tweet. If a tweet has a percentage of english words greater than a threshold σ we keep it, otherwise we drop it.

From the pruned tweets we extract the users, and we assign to the cleaned text the polarities found in the previous step.

2.2 Compute statistics

After having cleaned the raw data and after evaluating text polarities, we compute the statistics using the tuples we have kept in the aforementioned steps. First, we load and match the tuples with the tweet timestamp we get from the database. For each timestamp, we keep the date and time up to the hour, this means that any further analysis is condensed in a time window of 1 hour. Second, we match the hashtags with the polarities in order to understand positive and negative trends of the topics. It is important to notice that in a tweet more than one hashtag can be present, and we want to analyse each hashtag separately. This step is performed by a match operation with the sentiment polarities followed by a sequence of reduce operations to compute different statistics. The first information we want to extract is about the evolution of the popularity of an hashtag. To obtain this information we will count hour by hour how many tweets contain each hashtag, and similarly we also keep track of how many different users tweet about the same hashtag. In the same way we are keeping track, hour by hour, of how popularity changes for each hashtag. Then we collect for each hashtag the moment in which it reached his peak in popularity alongside the date of first and last appearance, marking in this way the lifespan of the hashtag. Sentiment analysis is performed not only to compute polarities of each tweet, but for each hashtag we aggregate three values, namely the average positive sentiment, the average negative sentiment and the emotional divergence present in each of them. These statistics can give us some insights on people's opinion towards the topic to which the hashtag is referring to. In particular, in [7] is demonstrated how the overall sentiment (polarity) in one tweet does not influence

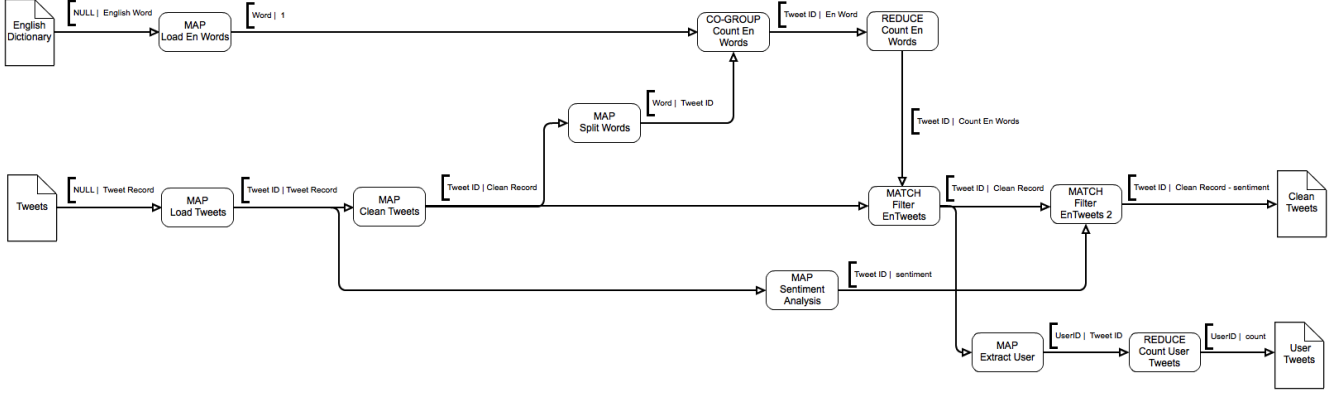


Figure 1: Data cleaning PACT

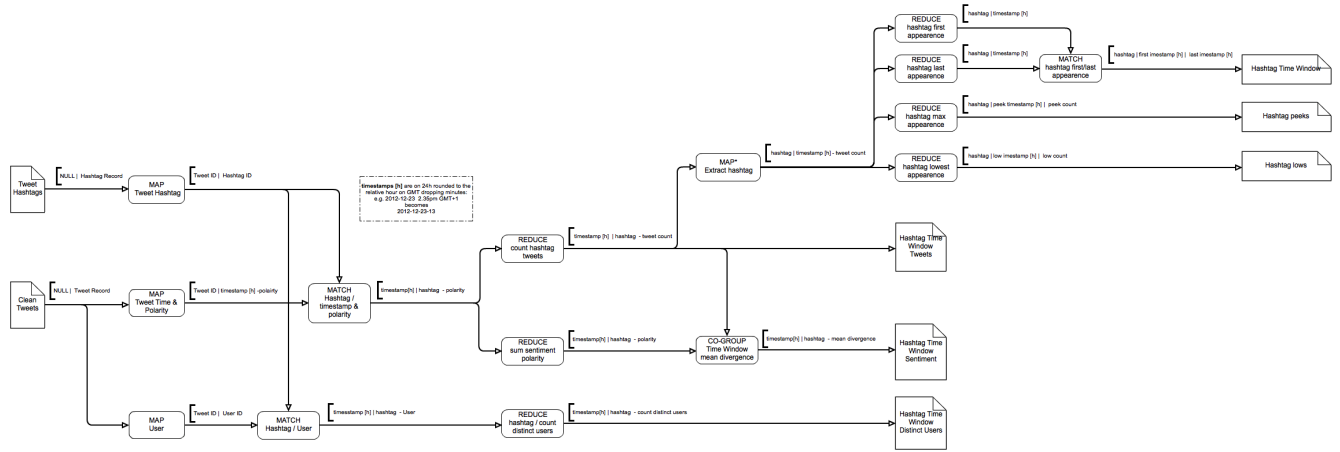


Figure 2: Compute statistics PACT

tweet id	user id	tweet
292375792485298176	858488612	#KidCudi - #EraseMe - The Whizz Bells : http://t.co/Lp1zABOV via @youtube
292375792481099777	486970282	RT @Kirra_: Today is a day where I need to crawl into my bed and sleep the day away.
292375792481099776	336390437	There are poor people, money is the only thing they got.
292375792476909568	74570186	I can't do anything without listening to music while I do it.

Table 1: A sample of the tweet dataset

its probability to be reshared. For this reason they propose emotional divergence as a new measure able to identify which tweets have a better chance of being retweeted. Aggregating this information for each hashtag we can gather additional informations in support of one hashtag being relevant and interesting.

3. DATA

We downloaded 33 million of English tweets from the tweet-stream using the streaming APIs⁵ on a period of approximately two weeks from December 18, 2012 to January 18, 2013⁶. A summary of the main characteristic of the dataset

⁵<https://dev.twitter.com/docs/streaming-apis>

⁶We didn't download any data during the week between the 25-12-2012 and the 2-1-2013

is represented in table 2. The number of users is 16 million, with an average of 2 tweets per user. The number of hashtags is one order of magnitude less than the number of tweets meaning that users often talk about same topics or do not mark their tweets with hashtags. From this dataset we generated three datasets 100K-TWEETS, 1M-TWEETS, 10M-TWEETS with 100 thousands, 1 million and 10 million tweets respectively, in order to test performance with various sizes of the dataset. We also downloaded a dictionary of 213 thousands English words to filter the meaningful tweets from those irrelevant or containing only symbols. To perform the sentiment analysis and (i.e., to assign a positive or negative polarity to the tweets) we used the SentiStrength library registering both positive and negative polarity carried in each tweet.

Period	2012-12-18/2013-01-18
Tweets	33774428
Users	16099129
Hashtags	1194691
Max tweets per user	2380

Table 2: Characteristics of the dataset used in the experiments

4. EXPERIMENTAL EVALUATION

We tested our solution with Stratosphere 0.2.1⁷ running on a GNU/Linux machine with 2Gb RAM DDR2, and a CPU AMD AthlonTM 64bit X2 Dual Core Processor 5000+⁸. We used a version of Stratosphere which is not publicly available at the time of writing as still under beta-testing. Bugs present in the implementation required us to have it running with a degree of parallelism set to 1.

All the experiments have been performed on samples of the original database, in order to show time and quality performance. We used Java JRE 1.6.0.26 to program the PACT using the Stratosphere libraries and server, additionally we integrated the SentiStrength java library⁹. We performed preliminary experiments to set the english word threshold and we finally set it to $\sigma = 0.1$, i.e. requiring that at least 1 word out of then would be a valid English word, in order not to prune too many tweets. Since we are matching with english words that are not stemmed a larger threshold removes interesting tweets.

4.1 Time performance

In this section we briefly report and comment time performance registered in the various experiments. We performed a total of 4 different experiments, where the flow and so the statistics retrieved were always the same, but we changed the set of tweet over which we computed them.

As said before we performed this experiment with a degree of parallelism se to 1. This means that the Stratosphere engine had been instructed to compute each step sequentially without exploiting possible parallel computations. In particular we had our flow running over the last 100 thousands, 1 million, 10 million and 20 million tweets, in reverse chronological order. Time results are present in table 3.

As expected, we registered time performance rapidly increasing with respect to the number of tweets analysed. Anyway, given the size of the dataset and the type of complex statistics computed it has in the worst case a total running time of at most 1 hour for 20 millions tweets.

4.2 Output charts

As described previously, the scope of this work is building a flow that can cope with an large amount of tweets, and which can parse, clean and filter them. Moreover this flow is designed to compute and aggregate statistics for the hashtags present in the flow and can be extended to many more

Number of Tweets	Time (secs)
100 thousands tweets	50
1 million tweets	183
10 million tweets	1516
20 million tweets	4006

Table 3: Time performances for the computation of the statistics with respect to the size of the dataset analysed

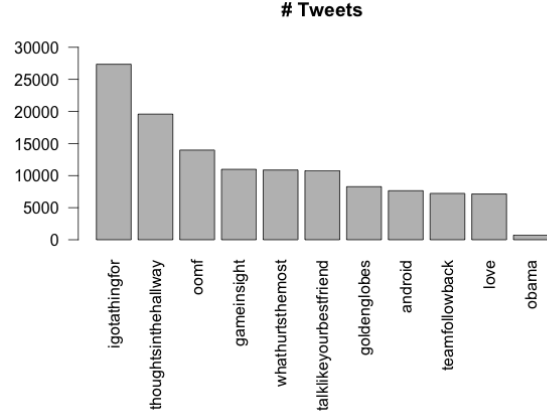


Figure 6: Top10 hashtags for popularities within tweets, plus the hashtag for “obama”; reported with the number of tweets in which they appear

analysis. Our initial assumption is that we can consider hashtags as indicators of events or discussion topics.

For instance we can consider the evolution of some hashtags, (in Figure 6 we show the Top10 hashtags for popularity, plus the hashtag for “obama” as a possibly interesting case) and we can analyse their lifespan as in Figure 7. Indeed identify-

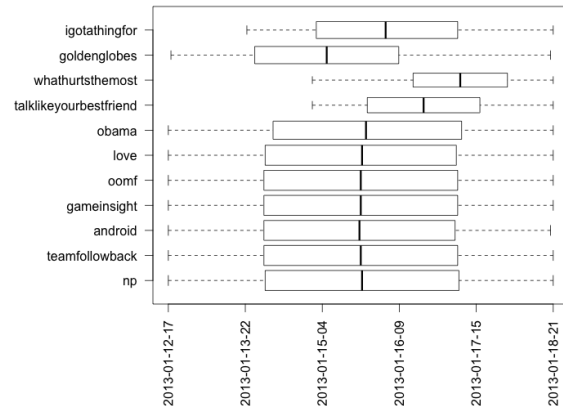


Figure 7: Lifespan of the Top10 hashtags plus “obama”, the thicker line delimits where the first and third quartile of tweets for that hashtag lies

⁷<https://stratosphere.eu/downloads>

⁸Linux 3.0.0-12-generic #20-Ubuntu SMP x86_64

⁹<http://sentistrength.wlv.ac.uk>

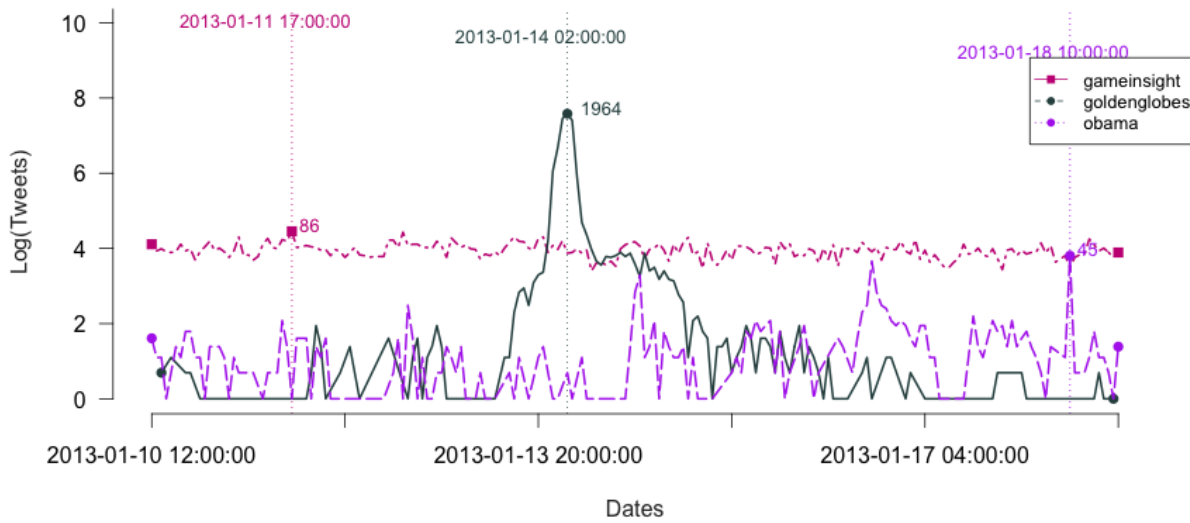


Figure 3: Popularity over time of three hashtags: gameinsight, goldenglobes, and obama. Number of tweets have been normalized through the \log function

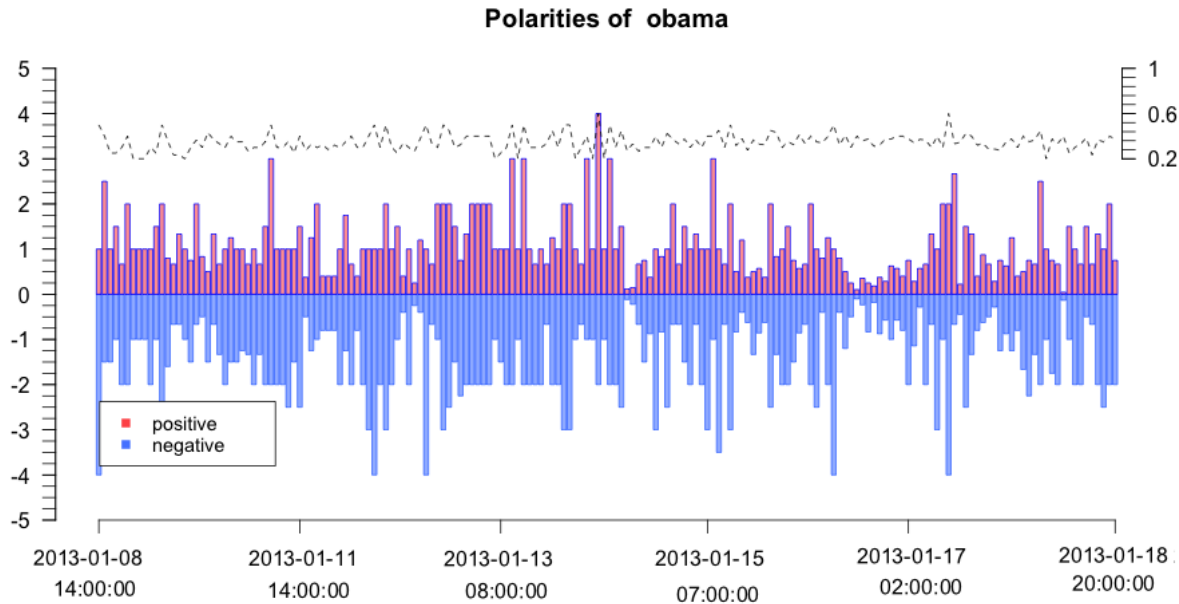


Figure 4: Sentiment polarity expressed through tweets containing the hashtag “obama”. The dashed line on the top, represent in each point in time the average emotional divergence.

ing that an hashtag has a limited lifespan within the flow of tweets may be evidence of an event happening in that particular time window. Looking instead to the actual evolution, hour-by-hour, of the popularity of each tweet we can identify interesting trends, points of peaks and lows. In Figure 3 we can see how the popularity of some hashtags (expressed as hourly number of tweets), can have very different evolutions. In the figure the hashtag “gameinsight” refers to the

name of one developer and publisher of mobile games and social games, they and their followers use this hashtag when referring to some of these games. We can see that during the time of the analysis the number of tweets mentioning this tag each hour is pretty stable. Instead, the other two hashtags have some highly unstable behaviour. The one referring to Barack Obama is not so much popular, even considering that there are more than one hashtag referring to President

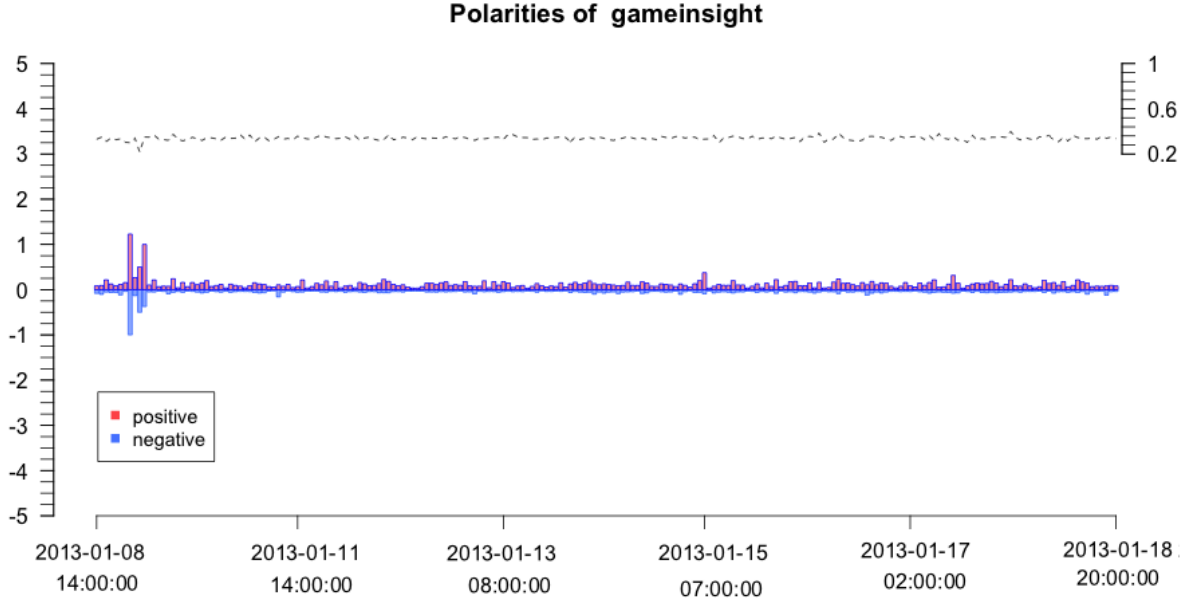


Figure 5: Sentiment polairty expressed through tweets containing the hashtag “gameinsight”. The dashed line on the top, represent in each point in time the average emotional divergence.

Hashtag	Number of Tweets
obama	1771
obamacare	173
nobama	120
impeachobama	73
obama2012	63

Table 4: Hashtags referring to “Barack Obama” and the number of tweets containing them

Obama, as we shown in Figure 4. In contrast to the previous hashtag, the hourly number of tweets containing the tag “obama” is more unstable, and we can see that there are time windows where this number is actually zero. Finally, in Figure 3, we analyse also tweets regarding the “goldenglobe” hashtag, referring to the annual Golden Globe Awards, held in 2013 on the night of 13Th. of January. The data-peek we can observe is indicative of when the award night where held, and we can easily see how before and after that peek was and returned to be pretty low.

Finally, another information to take into account, about the evolution of an event, is the polarity expressed in tweets about it. In Figure 4 and 5 we show average emotional polarities present in tweets in the hourly time-window. These pictures demonstrates how while tweets containing the “gameinsight” keyword are more popular, tweets tagged with the hashtag “obama” have a more strong weight in terms of sentiment polarities. These statistics combined with the evolution of popularity of hashtags are starting points for more complex models and analysis. They can be used to actually identify which hashtags in which time windows are likely to be indicator of events happening. Additionally, approaches like the ones presented in [10] can be applied not only to identify burst in popularity or in sentiment polarities, but

also to identify overlapping bursts, i.e. bursts in popularity for different hashtags in the same time window, so that we can actually try to detected which hashtags may refer to the same event.

In the end of our process we where able to aggregate, for the biggest dataset, statistics for more than 480 thousands hashtags, keeping tweets from more than 604 thousands users over 247 time windows. These aggregated statistics can easily be used as grounds for more complex and structured reasoning, but they are already able to highlight popular hashtags

5. ISSUES AND CONCLUSIONS

During the design and the main development of the project we didn’t encounter big difficulties or problems with the software provided. We found the PACT model intuitive and flexible for our needs, covering the design and the first implementation without any big issues. The software demonstrated fair performance even on the relative small machine at our disposal, coping with 20 millions tweets and being able to filter and join, for the more loaded component of the flow, more than 90million tuples.

Problems where encountered during the final testing of the application. Firstly we reported the following exception:

```
java.lang.RuntimeException:
Data distribution must not be null when
the ship strategy is range partitioning
```

which was stopping the Nephele runtime, and it web panel, from executing the PACT program submitted.

Later we noticed some non-coherent results being returned

as output of subsequent computations. In this second case, we found that the same PACT program, executed over the same machine with the same input, where actually computing different results in 1 out of 5 subsequent computations on average. Debug of this problem required time due to the current absence of a complete logging environment for the system, in its current version. After extensive analysis we traced back the origin of the problem to a mismatch in attributes used as field for an equijoin in the early steps of the flow, still we are not able to explain the apparent non-deterministic behaviour of the system.

6. REFERENCES

- [1] E. Bakshy, J. M. Hofman, W. A. Mason, and D. J. Watts. Everyone’s an influencer: quantifying influence on twitter. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 65–74. ACM, 2011.
- [2] L. Hong, A. Ahmed, S. Gurumurthy, A. J. Smola, and K. Tsioutsoulis. Discovering geographical topics in the twitter stream. In *Proceedings of the 21st international conference on World Wide Web, WWW ’12*, pages 769–778, New York, NY, USA, 2012. ACM.
- [3] J. Lehmann, B. Gonçalves, J. J. Ramasco, and C. Cattuto. Dynamical classes of collective attention in twitter. In *Proceedings of the 21st international conference on World Wide Web, WWW ’12*, pages 251–260, New York, NY, USA, 2012. ACM.
- [4] K. Lerman and R. Ghosh. Information contagion: an empirical study of the spread of news on digg and twitter social networks. *CoRR*, abs/1003.2664, 2010.
- [5] A. Marcus, M. S. Bernstein, O. Badar, D. R. Karger, S. Madden, and R. C. Miller. Twitinfo: aggregating and visualizing microblogs for event exploration. In *CHI*, pages 227–236, 2011.
- [6] M. Mathioudakis, N. Koudas, and P. Marbach. Early online identification of attention gathering items in social media. In *Proceedings of the third ACM international conference on Web search and data mining, WSDM ’10*, pages 301–310, New York, NY, USA, 2010. ACM.
- [7] R. Pfitzner, A. Garas, and F. Schweitzer. Emotional divergence influences information spreading in twitter. In *ICWSM*, 2012.
- [8] M. Thelwall, K. Buckley, G. Paltoglou, D. Cai, and A. Kappas. Sentiment in short strength detection informal text. *J. Am. Soc. Inf. Sci. Technol.*, 61(12):2544–2558, Dec. 2010.
- [9] A. Tumasjan, T. O. Sprenger, P. G. Sandner, and I. M. Weppe. Predicting elections with twitter: What 140 characters reveal about political sentiment. In *Proceedings of the fourth international aaai conference on weblogs and social media*, pages 178–185, 2010.
- [10] M. Vlachos, C. Meek, Z. Vagena, and D. Gunopulos. Identifying similarities, periodicities and bursts for online search queries. In *Proceedings of the 2004 ACM SIGMOD international conference on Management of data, SIGMOD ’04*, pages 131–142, New York, NY, USA, 2004. ACM.