

# Spatio-temporal analysis of opinion in social media: outlier detection for the business intelligence area

Miguel Fernandes Caíña, PhD Student

Supervised by: Dr. Rebeca P. Díaz Redondo and Dr. Ana Fernández Vilas

Department of Telematics Engineering, University of Vigo



## Motivation

- Internet has become a communication and expression platform, rather than just a static information source. Mailing lists, forums and chats have been part of it since the beginning, but over the last years, social networks have become the primary platform of communication for the majority of its users.
- The continuous flow of public information from forums and social networks makes possible to extract any kind of sentiment expressed about a product, service or brand. The aggregation of this data, including the impact of time and location, could be crucial in the success of a business decision.
- Natural Language Processing (NLP), defined as the ability of a system to process human language (Preeti & BrahmaleenKaurSidhu 2013), is an artificial intelligence component that can be used to mine opinion and sentiment from social networks, and classify each post as being positive, negative or neutral towards a specific subject.
- This flow of opinions could be exploited by a Company, in order to verify the impact of a business decision or an external situation on the public's perception over its products, services or brands. Simply ignoring it could be harmful for the company's success.



## Objectives

- The main objective of this PhD is to propose a general model applying different techniques (opinion mining, relevant topic identification, data and company connections, space-time scopes and real time analysis) that could be transformed into a solution that provides a high level view about products and services of interest, enabling decision making "as soon as possible", as well as post-mortem analysis of relevant events.
- A platform will be developed following the proposed model, and it should be able to provide insight about the impact of events on the public's perception over related brands, services or products.

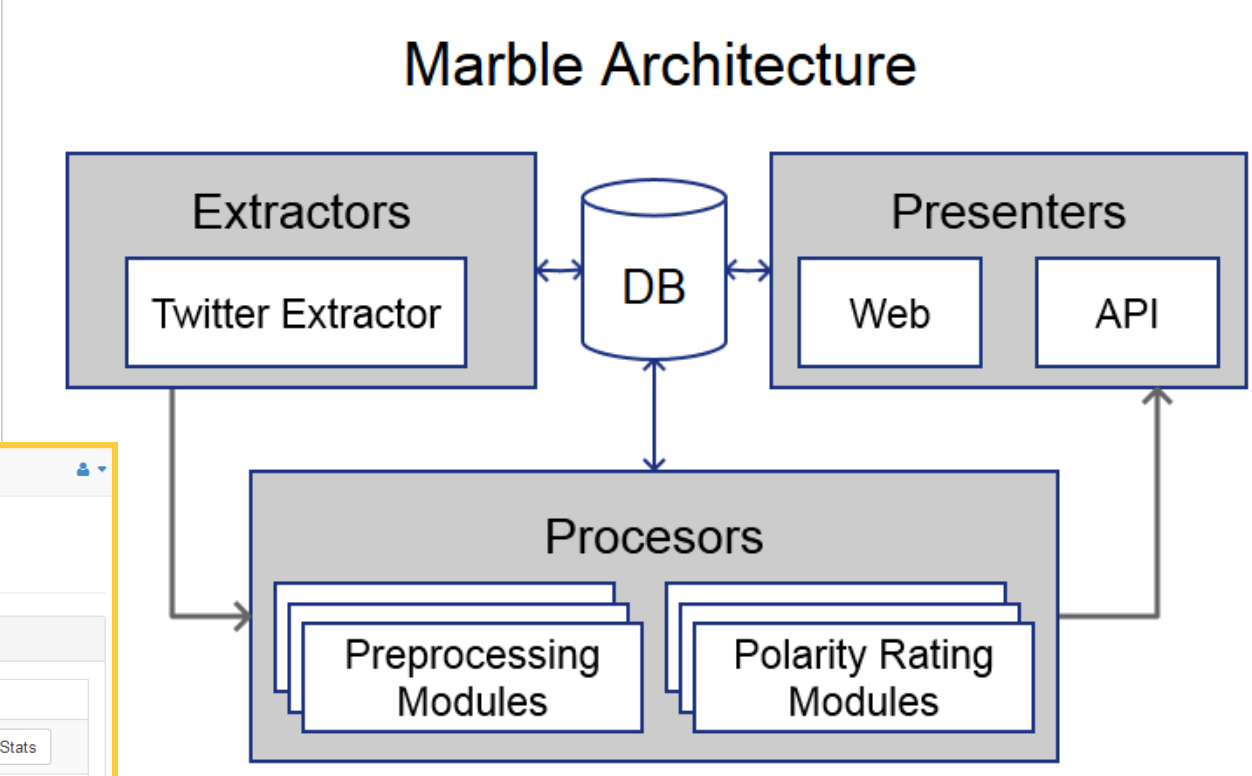


## Research Plan

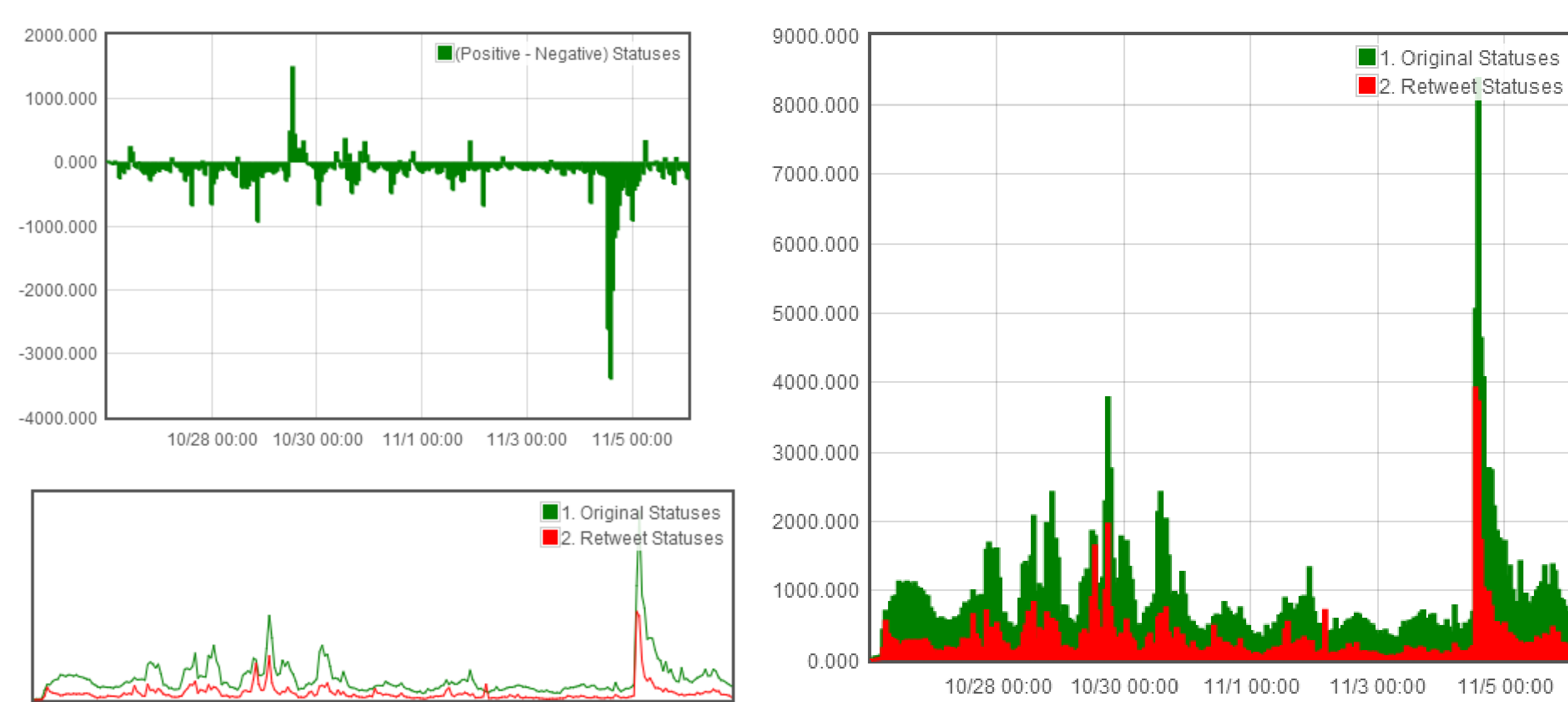
- A comprehensive review of the relevant literature in the fields of opinion and sentiment mining, topic disambiguation, space-time scopes and real time analysis, outlining the state-of-the-art in the fields is being performed in order to gain critical insights.
- A platform, named **Marble** has been defined, and integrates all the main components that could provide the achievement of the objectives. This model is flexible enough to allow an implementation based on its guidelines. The main components are:
  - + Subject selection and data extraction: an extraction module capable of extracting information from Twitter related to a subject. If needed, the module will be expanded to cover other social networks.
  - + Data preprocessor: a module incorporating NLP processing techniques, stemming and lemmatization capabilities, synonyms recognition and disambiguation practices, that will be in charge of converting the raw data into information for the opinion mining module. This module would be configurable, in order to use different processing techniques, depending on the nature of the data.
  - + Sentiment Analysis: a module in charge of extracting the opinion expressed in each message and define a polarity, using a combination of sentiment analysis techniques and heuristics, which will allow to identify specific characteristics of opinion on each user interaction.
  - + Presenter: a presentation module, responsible for extracting relevant information from the mined opinion and correlating it to manually identified events. The module will also be able to detect "special situations" not related to any of the known events, in order to discover unidentified incidents.
- The development is following an iterative and incremental approach, and is divided in several phases:
  - + Phase 1 - Conceptual Test: A basic approach using simple scripts written in Perl, to extract user's tweets, preprocess them and mine their opinions using a basic algorithm.
  - + Phase 2 - Platform Definition: Definition a Java enterprise application model, using MongoDB and PostgreSQL as data warehouses, and an extensible architecture to be able to accommodate the four main modules, and their subcomponents.
  - + Phase 3 - Basic Modules: Development of basic algorithms for each components, following the same principles as in Phase 1.
  - + Phase 4 - Advanced Modules: Integration of several state-of-the-art techniques for each area, specially in preprocessing and sentiment analysis. Automatic anomalies detection, external API support and live extraction will also be part of this phase.
- Validation and assessment of the results will be based on a statistical approach, and the project success will be evaluated using this criteria.

## Preliminary Results

- Phases 1 to 3 have been completed.
- A conceptual test have been performed using shell scripts.
- The **Marble** platform have been implemented and tested with the data extracted in the conceptual test.
- Two subjects were selected for the experiment: *Blackberry* and *Whatsapp*.
- Opinion and impact of events have been extracted, and the results of the experiment were submitted to the KDIR 2014 (Scopus Index) conference for its approval.
- The development is being hosted on Github, and the source code is publicly available at: <http://iclab.det.uvigo.es/marbleproject.html>

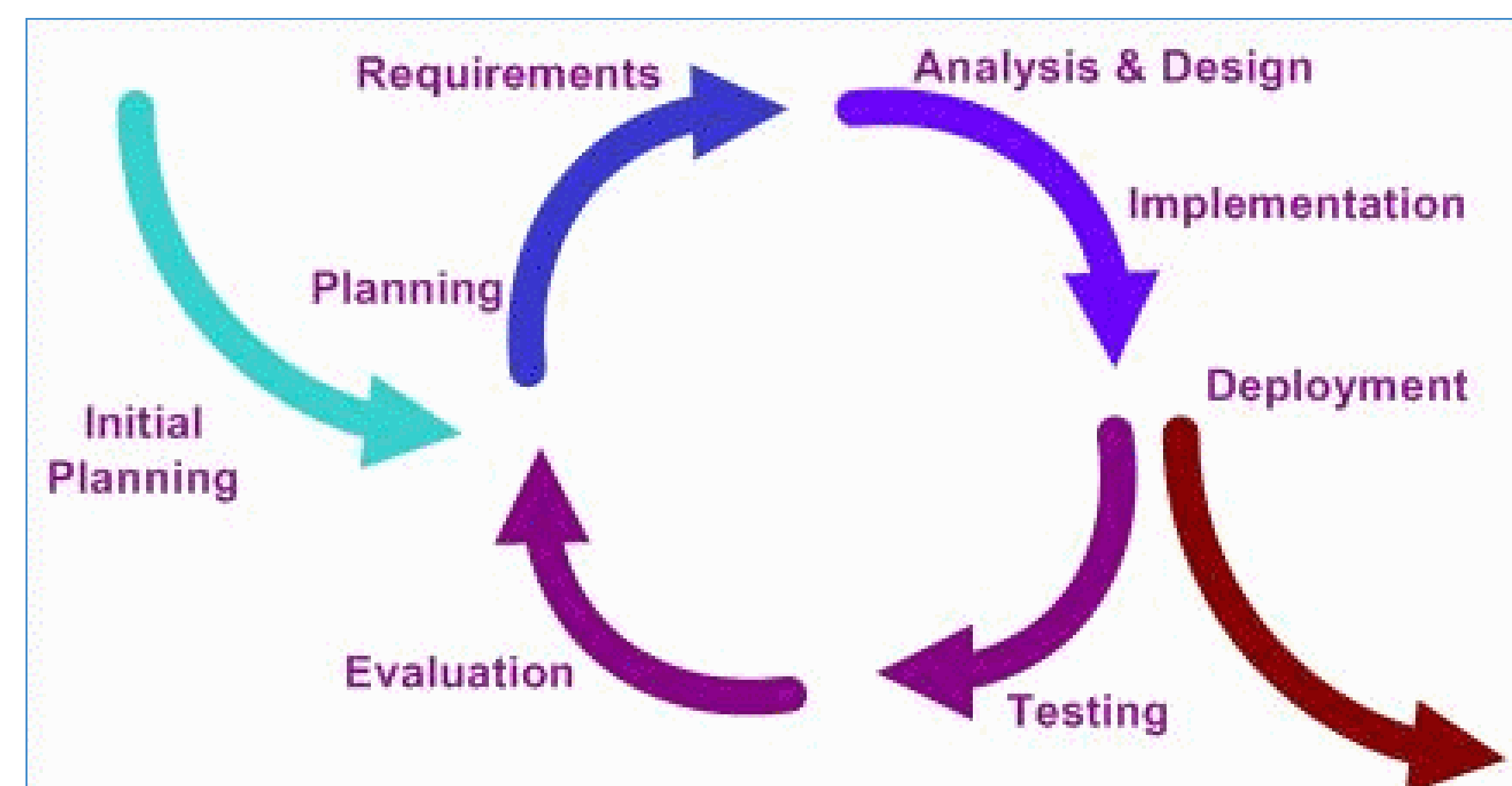


Time	Statuses	Originals	Retweets	Positive Statuses	Negative Statuses	Positive - Negative	Average Polarity	Unique Users
Wed Nov 06 2013 02:00:00	707	601	106	216	472	-256	0.2027004350545917	703
Wed Nov 06 2013 01:00:00	864	675	189	353	492	-139	0.22643301644919026	841
Wed Nov 06 2013 00:00:00	983	755	228	450	508	-58	0.2560683200601234	942
Tue Nov 05 2013 23:00:00	1085	849	236	480	565	-85	0.23172685374267324	1030



## Next Year Planning

- The results obtained until Phase 3 will be tested and evaluated using statistical methods on real data extracted from social networks, and the analysis will be submitted to a journal (JCR indexed).
- Start of Phase 4 development, including automatic anomalies detection and a selection of advanced preprocessing and opinion mining techniques.
- Improvement of the platform and publication of internal documentation.
- Attendance to at least one international conference or workshop.



Iterative and Incremental Development

## References

Gokulakrishnan, B. et al. 2012, "Opinion mining and sentiment analysis on a Twitter data stream", *Advances in ICT for Emerging Regions (ICTer)*, 2012 International Conference on, Dec, pp. 182.  
 Cambria, E., Havasi, C. & Hussain, A. 2012, "SenticNet 2: A Semantic and Affective Resource for Opinion Mining and Sentiment Analysis.", *FLAIRS Conference*, eds. G.M. Youngblood & P.M. McCarthy, AAAI Press.