# Sentiment Analysis on Social Media

Dr.K.Dharmarajan[#1], Farhanah Abuthaheer[*2], Dr.K.Abirami[@3]

[#1]Department of IT, VISTAS, Chennai, Indian

[*21]Department of IT, VISTAS, Chennai, Indian

[@3]Department of BCA, Guru Nanak College, Chennai, Indian,

*Abstract*— **Microblogging platforms like Twitter can convey short messages to direct contacts, but  also to other potentially interested users.  They are actively exploited either by individual users or whole organizations and companies. This paper describes some results we obtained from the Social Network and Sentiment Analysis of a Twitter channel, related to a pop music event. Apart from the particular results a methodology and some guidelines for the automatic classification of Twitter content are discussed.**

*Keywords*—— **Social Network, Sentiment Analysis, Hierarchical Classification**

## I. INTRODUCTION

In the common meaning of the term, an online community (or virtual community) is a  group of people interested  in  a particular topic, or that share some ways of thinking, or that in general have some kind of link that brings them together, with the peculiarity that they interface  and connect to  each other through data communication network (such as Internet). In this  way, they form a  social  network with unique characteristics:

Fact this combination  is  not necessarily bound to a physical place and anyone can participate wherever he is, with a simple access to networks. The social networking sites (SNSs), as  defined by Boyd and Ellison in [1], are a collection of web-based services that allow users to build a profile within the system and define a list of other users with whom they have some kind of connection. According to Sunden profiles are unique pages where one can "type oneself into Being" [2], as the creation of a profile is the minimum condition for joining a SNSs. What makes the SNSs unique is that their purpose is not in most cases, to allow users to make new  friends but the emphasis  is on making visible their existing social networks and on  the chance to describe them.  On the other hand, the specific features of each social network site may depend  also on the  possible target (social, linguistic or geographic) to which the service is directed. The architecture  of social  networking platforms is  much differentiated. While the most popular platforms are built  as  essentially centralized systems, other platforms have  a distributed architecture [4].

The decentralized systems, in particular, often use some notion of trust and cryptography to address the risks of online social networks, which are perceived as serious by many users and have led to incidents.

Ethnicity, religion, sexual orientation, political beliefs are other  factors that have  led to the establishment  of dedicated social network services, but probably they are also playing an active   role in creating and aggregating online communities leveraging the   bigger and most popular social networks. This suggests the possibility of new ways to  spread information  and  to influence public opinion [2] [3]. These new scenarios can be better evaluated by a combined observation of the structure and the actual content  of the network.  This  kind of  analysis could   highlight  emerging  social behaviours. for example,  the possible  differences in the  sentiment polarity of  female and male users, towards the discussed topic are examined.

To investigate on  the content  and  on the relations among the actors of a network, it could be useful to  contextualize the  network  itself.  In particular, it could  be important to  consider and inquiry  the  content  of  the  messages  that guide the relationships of the community. It is only through this kind of investigation that we can analyze the semantic meaning of a link, from which we  could  infer  the  kind  of  relationship.  This sharpens our description of the social network in many of its facets.  A useful tool for such  surveys is  Sentiment Analysis (SA). SA is  a branch of Opinion Mining that aims to listen and process the data that users post on social media.  It is    an interdisciplinary field that in recent years has had a significant growth and that makes extensive use of machine learning techniques [5].

The information about social relationships can be used to improve user-level sentiment  analysis. Sentiment Analysis  is mapped on  social media with observations    and    measurable data;   the results highlight the importance of SNSs (i.e. Facebook) as a platform for online marketing. Social media provide tremendous challenges for researchers and analysts trying to gain insight into

human and group dynamics. A central problem is the sheer amount of data available in social media. For example, the social media aggregation site Spinn3r [3] advertises that they provide information on over 20 million blogs feeds that stream over 100 thousand posts and an 8 month archive consisting of 21 TB of data. Facebook currently involves over 400 million active users with an average of 120 `friendship connections each and sharing 5 billion references to items each month. One analysis approach treats the interactions as graphs and applies tools from graph theory, social network analysis, and scale-free networks. However, the volume of data that must be processed to apply these techniques overwhelms current computational capabilities [9].

Even well-understood analytic methodologies require advances in both hardware and software to process the growing corpus of social media. Social media provides staggering amounts of data [10]. Extracting knowledge from these volumes requires automation. Computing quickly over this data is a challenge for both algorithms and architectures.

## II.  SENTIMENT ANALYSIS ON TWITTER

In this research work, we built a system for social network and   sentiment analysis, which can operate  on Twitter data. Twitter is a popular Platform for social networking and micro blogging, counting hundreds of millions of active users and daily published messages. As a social networking platform, Twitter is structured as a directed graph, in which each user can choose to follow a number of other users (followers), and can be similarly followed by other users (followers) [11]. Thus, the "follow" relationship is asymmetrical, it does not require mandatory acknowledgement, and it is essentially used to receive all public messages published by a followed user. As a micro blogging service, Twitter is used to publish short messages counting a maximum of 140 characters (tweets), which may contain opinions, thoughts, facts, references to images and other media. Moreover, through the @ symbol it is possible to introduce mentions i.e. references to other users, and through the # symbol it is possible to introduce hash tags, i.e. references to discussion topics [9]. Sentiment Analysis can be done at document, phrase and sentence level. In document level, summary of the entire document is taken first and then it is analyse whether the sentiment is positive, negative or neutral. In phrase level, analysis of phrases in a sentence is taken in

account to check the polarity. In Sentence level, each sentence is classified in particular class to provide the sentiment. Sentimental Analysis has various applications. It is used to generate opinions for people of social media by analysing their feelings or thoughts which they provide in form of text. Sentiment Analysis is domain centred, i.e. results of one domain cannot be applied to other domain [12].

Sentimental Analysis is used in many real life scenarios, to get reviews about any product or movies, to get the financial report of any company, for predictions marketing. Twitter is a micro blogging platform where anyone can read or write short form of message which is called tweets. The amount of data accumulated on twitter is very huge. This data is unstructured and written in natural language [8] [9]. Twitter Sentimental Analysis is the process of accessing tweets for a particular topic and predicts the sentiment of these tweets as positive, negative or neutral.

Consequently, in our analysis we collected  three types of data. The User type represents users' profiles; from Twitter we obtain the following fields: user_id  ,name,  location,  num_  followers, num_tweets. The Tweet type represents posted messages; Twitter  we  obtain the following fields: tweet_id, user_id,  message, date. Finally, the Friend  type represents the "follow" relationships among  users. Apart from data  obtained directly from  Twitter, we added a  field to both tweetsand users, to associate a sentiment with them, according to the result of our analysis.

## III. PROPOSED WORK

### A.  System Analysis

The main objective of this thesis work is to perform the sentiment analysis on Indian Political Parties like BJP, INC and AAP, such that people opinions about these parties progress, workers, policies, etc. which are extracted from Twitter. Thus to achieve this objective we build a classifier based on supervised learning and perform live sentiment analysis on data collected of different political parties To achieve this objective the following methodology is used: A thorough study of existing approaches and techniques in field of sentiment analysis. Collection of related data from Twitter with the help of Twitter API Pre-processing of data collected from Twitter so that it can be fit for mining.

- Computing the result of different classifier using dataset collected from Twitter.
- Comparing results of each classifier and plotting a graph that show the trend of positive and negative sentiment for different political parties.

The analysis of political communication in the mass media. This study takes also in account the data provided by Sentiment Analysis is a process of extracting feature from user's thoughts, views, feelings and opinions which they post on any social network websites. The result of sentiment analysis is classification of natural language text into classes such as positive, negative and natural. the analysis of Social Media, of Facebook in particular, with measurable data, available to public domain.

An ability to classify tweets by location in real time is crucial for applications exploiting social media updates as social sensors that enable tracking topics and learning about location-specific trending topics, emerging events and breaking news. Specific applications of a real-time, country-level tweet geo location system include country-specific trending topic detection or tracking sentiment towards a topic broken down by country. To the best of our knowledge, our work is the first to deal with global tweets in any language, using only those features present within the content of a tweet and its associated metadata.

### B.  Collaborative Filtering Algorithm:

If the user is giving any query, for that query all the tweets are taken from the live twitter. Then by using the collaborative filtering algorithm, we classified the tweets into Happy, Sad and Funny. If the person tries to tweet bad words means, the system cannot allow the user to post in twitter.

A Classification Engine which classifies search results into clusters and sub-clusters recursively, highlighting meaningful relationships among them, or assigns documents to predefined thematic groups.

- Easy to monitoring the social media
- Interact with mining classify to social media Happy, Sad and Funny Tweets are predicted.

### IV. PROPOSED MODULE DESCRIPTION

### A.  User Profiling:

User profiling is the foremost part of our project which covers the single sign on user and also registers new user to the ecommerce site.
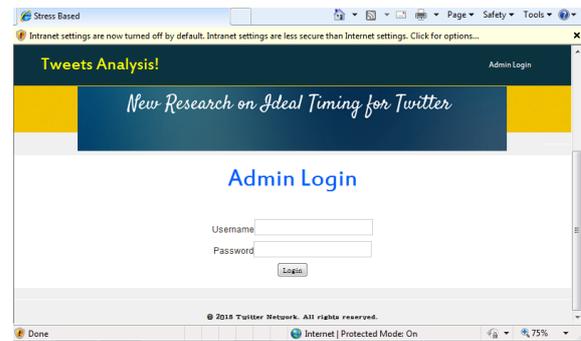


Fig 1.Twetter Analysis Login Page

### B.  Pre processing

In this module, we are collecting tweets dataset from live twitter by creating the app in apps.twitter.com . The dataset which contains tweets, user id, time/date and username which was given by the user are imported into the database. Then it is processed by stop words and stemming technique. The pre-processing phase The Fig 2. Shows the Load Twitter Dataset



Fig 2.Load the Twitter Dataset

### C. Stop words Removal:

A dictionary based approach is been utilized to remove stop words from tweets. A generic stop word list containing 75 stop words created using hybrid approach is used. The algorithm is implemented as below given steps. The target text is tokenized and individual words are stored in array. A single stop word is read from stop word list. The stop word is compared to target text in form of array using sequential search technique. If it matches , the word in array is removed , and the comparison is continued till length of array. After removal of stop word completely, another stop word is read from stop word list and again algorithm runs continuously until all the stop words are compared .Resultant text devoid of stop words is displayed, also required statistics like stop word removed, no. of stop words removed

from target text, total count of words in target text, count of words in resultant text, individual stop word count found in target text is displayed. The Fig 3 shows Pre Processed Twitter data



Fig 3. Preprosed Twitter data

### D. Stemming Technique:

After removing the unwanted words from the tweet, stemming technique is processed. Stemming is the process of reducing inflected (or sometimes derived) words to their word stem, base or root form generally a written word form. The stem need not be identical to the morphological root of the word; it is usually sufficient that related words map to the same stem, even if this stem is not in itself a valid root.

### E.  Classification:

After stemming process, all the reviews terms are classified into happy, sad and funny tweets. Collaborative filtering algorithm is used for classification. Here we are having happy words, sad words and funny datasets. By comparing with this, we can classify the tweets into Happy, Sad and Funny Tweets.
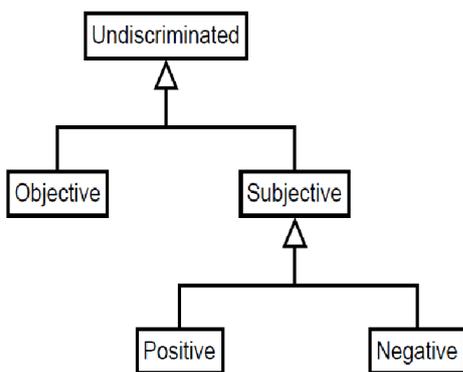


Fig 4. Classification of Sentiment Tweet

### F.  Capturing user Sentiments:

We present the main framework of our method. We regard extracting opinion targets/words as a co-ranking process.

We assume that all nouns/noun phrases in sentences are opinion target candidates, and all adjectives/verbs are regarded as potential opinion words, which are widely adopted by previous methods. Each candidate will be assigned a confidence, and candidates with higher confidence than a threshold are extracted as the opinion targets or opinion words.

If a word is likely to be an opinion word, the nouns/ noun phrases with which that word has a modified relation will have higher confidence as opinion targets.  If a noun/noun phrase is an opinion target, the word that modifies it will be highly likely to be an opinion word. It can see that the confidence of a candidate (opinion target or opinion word) is collectively determined by its neighbours according to the opinion associations among them. A noun/noun phrase can find its modifier through word alignment. The Fig 5. Shows the User opinion targets Sentiments Tweet.
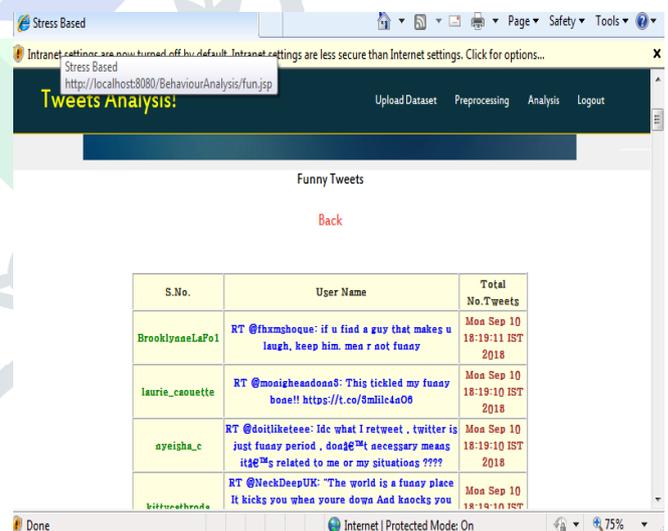


Fig 5. User opinion targets Sentiments Tweet

We additionally employ a partially-supervised word alignment model, which performs word alignment in a partially supervised framework.

### G.  Mapping Opinion Targets:

We finally map our feedbacks on positive, negative and neutral scale and then visualize this feedback in a graphical way.

As a communication medium, tweets have acquired peculiar nature. Some distinguishing features of communication  on Twitter are related to technical aspects; those include length of text, tags, urls, etc.. Other features may be classified as idiomatic use of the medium, and create a  sort of Twitter culture; those features include typical content and  most discussed topics, idiomatic expressions, abbreviated forms, etc. For example, a tweet may have the following form:

«RT  @richman  wow this  is  the #happiest day of  my life. #happy #glad #icantbelievit :) :D http://t.co/4VEH827bG7»

The peculiar nature of tweets requires specialized analysis techniques. As a  start, a tweet may contain  many  elements which are not significant for our classification, and can thus be dropped though a filtering process. To polish the message, we defined various filters, which we have applied  in  a customizable sequence. Proc. of the 16th Workshop "From Object to Agents" (WOA15) June 17-19, Naples, Italy54

A first filter eliminates  useless  tokens. Removed  tokens include:  the  starting  "RT" sequence,  which  indicates  a republished messages from a different user  (i.e. a retweet); the @ character and the whole following user name; the # symbol,  but not the following topic name, which is kept in the message. The topic name is  also removed,  though,  when it coincides with the name of the channel where  tweets are collected from. A second  filter  applies  the  language  specific rules.  It includes an orthographic correction of the message,  which is used to remove unknown words, which may not appear in any other tweet (in the example: "i cant believe it"). Ideally, the filter at this  level should also support  stemming  and removal  of stop words. However, those operations can be easily performed by  Weka, which we used for analysis. Finally, another filter  separates all punctuation  symbols from the text, and organizes them as  single-character words. However, some typical  patterns are  kept  as  aggregates, including smiles sequences, repeated question and exclamation marks. The final  result  of the filtering process  is  a  word vector, which  is then submitted to  the  classifier agents. As we  have mentioned, our  analysis aims at  identifying  the  following

classes  of messages:  un discriminated, objective, subjective, positive, negative.



Fig 6. User sad Tweet Classification

The system is organized as a   simple hierarchy of agents, mimicking the hierarchy  of  sentiment classes. In  fact,   since objective messages have no polarity by definition, the classifier for   positive and  negative  sentiments  is only  applied to subjective messages. If a  message  fails to be classified at the first stage, then it simply remains discriminated. If it fails to be classified  at the second  stage,  then  it  is  marked  as generically subjective.

Currently,   the   classifier agents apply   the Multinomial Naive Bayes algorithm, but other methods can be used and different agents  can be plugged in  the  system. However, instead  of generating a  training  set  by hand, we aimed at realizing an  automated (or  at  least  semi utomated) process for  obtaining good training sets. About  the objectivity/subjectivity classifier, we adopted  a similar strategy to [6]. In fact, to obtain objective content, we gathered messages generated from  popular news agencies. In our  tests,  we used  the  following  list: @ABC, @BBCNews, @BBCSport,   @business,   @BW,   @cnnbrk, @CNNMoney,   @fox32news,   @latimes, @nytimes,   @TIME.   To   obtain subjective content,  instead,  we gathered comments directed to the same list of users.

About  the  polarity  classifier, we  decided  to search  for sources of  mostly positive or negative

messages,  respectively. On the one hand, those sources should fit the particular setting of Twitter (short messages, idiomatic expressions, smiles, etc.) On the other hand, they should  not  be  specific to a particular topic or context (sport, music, etc.). Thus, we dropped the  idea of  collecting  messages about  particular events,  mostly generating either positive or  negative  sentiments. Instead, we collected messages, using generic yet polar terms as  queried hash tags.

In particular, we  used  the  following channels to gather positive content: #adorable, #awesome, #beautiful, #beauty, #cool, #excellent,  #great. We  used  the  following channels to  gather negative  content: #angry, #awful, #bad, #corrupt, #pathetic, #sadness, #shame. Actually, such terms have been chosen quite empirically, taking into account the quality  of  training  sets they generated.  But  they  could  be  selected  from WordNet-Affect  [8],  SentiWordNet [9],  and other affective lexicons, in a more systematic way.

This way, the training  set  is generated  in an automated fashion, as a list of tweets. Each tweet is associated with its supposed class, in accordance to its source. In fact, the training set is not perfect, as it contains messages gathered from public channels. However, a training set of this kind can be generated easily and in a methodical way, from real and updated Twitter messages. In the next section, we will also discuss the quality of results that can be obtained, using it as a basis for sentiment analysis. The Fig 7 Shows the User Classified  output report.
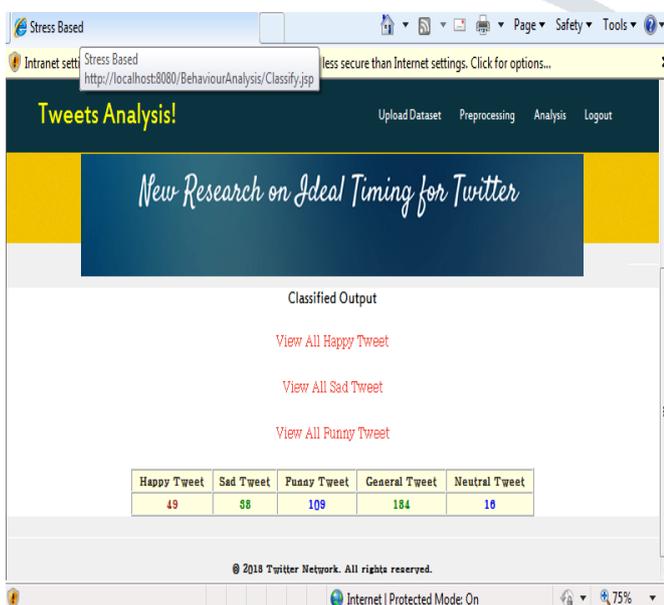
Fig 7. User Classified   output report

The  training set  can  be provided directly to the classifier agents. In the present form, the system is based on Weka, and can thus be configured for performing additional preprocessing steps  on  the messages,  including  common  TF-IDF transformations,  stemming,  elimination  of stopwords, exclusion of infrequent words, etc.

Currently, we  analyze tweets  for discriminating the basic classes of objectivity and polarity, at two levels. However, we designed  the  system  for more  complex  hierarchical classification,  with the  application of various types of classifiers, as an alternative to current Naive Bayes.

In fact,  hierarchical  classification  has  been applied successfully in a number of studies, for information retrieval. It  has  been  proven effective  especially  in  the  case  of classification  over  hierarchical taxonomies. Moreover, it has the advantage of being modular and customizable, with respect to  the  classifiers used at different levels. Using  the  same  probabilistic classifier  and  a  maximum  likelihood estimator, instead,  does  not  provide  advantages  for  the hierarchical  approach  over  the  flat  approach. Mitchell  has  proved  that  the  same   feature sets represent  documents  in  both  approaches. Consequently,  the  whole  hierarchical  classifier system   is  equivalent  to  the  corresponding  flat system.

Also  in  the  case  of  sentiment  analysis,  a hierarchy  of  classes  can  be  defined . Accordingly,  hierarchical  classification  has already been  applied  to  sentiment  analysis.
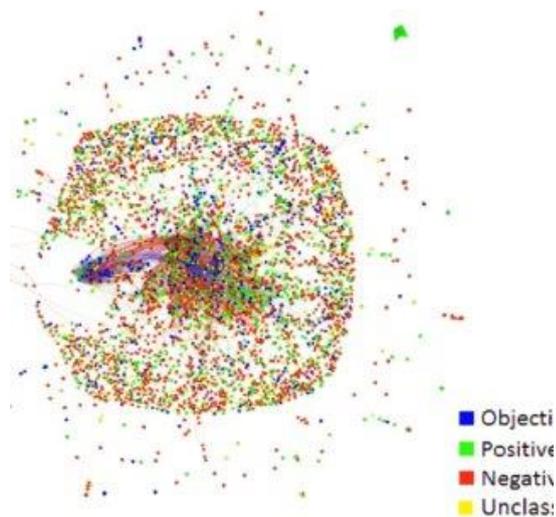
strong similarities found between the type of the published tweets and the instances used for training the classifiers. All data were downloaded between 2015-02-02 and 2015-02-10. The awarding of the Grammy took place on 2015-02-08. The network consists of a total of 5570 nodes and 6886 arcs.



Fig 8. User Classified report

### V.  CASE-STUDY: THE#SAMSMITH CHANNEL

This section will show the results of the classifiers and the analysis carried out on a case study.

With the  above described software, it is possible to obtain some  training sets  for  the classifiers. In our case  study, they consist of:

- 86000 instances (polarity)
- 32000 instances (subjectivity)

These instances have been obtained by exploring more than 60 channels on the social network.

In the generated models, the selected features are consistent with   our  expectations:  the  typical expressions   of   a   certain   feeling   (such   as smileys,     or     some     words     that     express appreciation or disgust) show a higher probability of belonging  to   the   class  of   that  feeling,   rather than  to  the  class  of the  opposite sentiment.

The   obtained   results   by   the   classifiers using   cross-validation   (with   folds   =   10)   on the  training sets showed  an accuracy of:

- 77,45% (polarity classifier)
- 79,50% (subjectivity classifier)

These results Fig 9. show that the model of the classifiers  contains  effective  features   for     the recognition  of  the  sentiment  of  a message..The case study  which  was  considered  in this  work  is the  social  network of  the #SamSmith channel (the singer who won four awards at the Grammy Awards 2015). The choice of this channel is justified by the



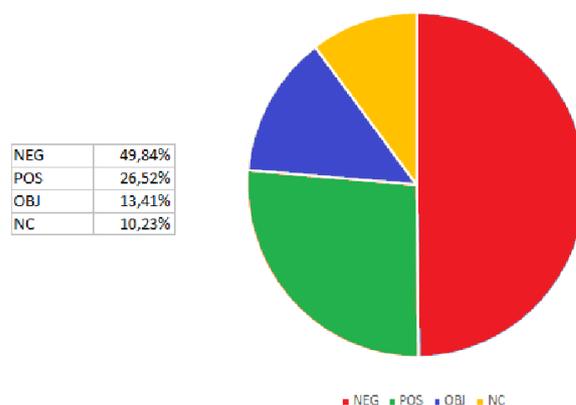| | |
|---|---|
| NEG | 49,84% |
| POS | 26,52% |
| OBJ | 13,41% |
| NC | 10,23% |

Fig 9. Graphical representation of recognition    of   the sentiment  of  a message

It    is    possible    to    notice    that    the    network topology  is  consistent  with  the  nature  of  the considered  case.  In  fact,  most  of  the  channel consists   of independent users (or small groups of users) that express their opinion about   the   artist; however,  in  the  central  part  of  the  network there are  some  major  communities.

The  prevailing  sentiment  detected  from  the classifier  is  the    negative  one.  Performing  an analysis  on  a  Communities  participating  in  the #SamSmith  channel.  Sentiment  analysis  on  the #SamSmith channel.

This  work  noticed  that  many  sentences  are actually  quotes  of  songs.  These  messages  contain melancholic  and    sad    phrases,  and  are  therefore classified as negative. Considering that a quote is generally  an  appreciation  for  the    artist,    most users  classified  as  negative  are  actually positive users. This is a typical example of a classic problem  of  misunderstanding  of    the  SA:  the system,  while  classifying  correctly    the  tweet, misses  the  assessment  of  the  feeling  because  it cannot evaluate the tweet together with its context. For   evaluating   the   performances   of   our system,    we  conducted  a  simple  survey  through  a group  of  persons  in  our  department.  In  this  way, we     selected     and     classified     100  messages  that show  a   clear  opinion  on  the   singer.  Then,  we

used those messages as a test. The results of the classifiers showed an accuracy of 84% for the polarity and 88% for subjectivity.

In the network periphery, it is possible to notice a small group of users whose feeling is completely positive. After a careful analysis of users' tweets in this small group, it was found that these posts are mainly retweets and the original messages are only two. Of these two messages, the first is actually positive, while the other one is objective. This episode shows how some errors of assessment can have important impact on larger communities.

## VI. CONCLUSIONS

In this article, we describe some results obtained from the synthesis of Social Network Analysis and Sentiment Analysis applied to the channel #SamSmith during the Grammy Awards in 2015. Apart from the particular results, a methodology and some guidelines for the automatic classification of Twitter content have been discussed.

The implemented software allows: (i) to get a training set for the classifiers that deal with Sentiment Analysis, and (ii) to make a thorough study of the topology of the networks. The study of the global sentiment within the network has highlighted the typical problems of Sentiment Analysis (irony, sarcasm, lack of information, etc.). Additionally, some peculiar problems of the considered channel were also detected (such as the quotes of songs).

The performances obtained by the classifiers during tests conducted on the training set and the analysis of the case studies have shown good and promising results.

## REFERENCES

[1] Boyd, Danah M., and Nicole B. Ellison. "Social network sites: Definition, history, and scholarship." Journal of computer‐mediated Communication 13.1 (2007): 210-230.

[2] Fornacciari, Paolo, Monica Mordonini, and Michele Tomaiuolo. "A case-study for sentiment analysis on twitter." WOA. 2015.

[3] Stefanone, Michael A., Derek Lackaff, and Devan Rosen. "The relationship between traditional mass media and "social media": Reality television as a model for social network site behavior." Journal of Broadcasting & Electronic Media 54.3 (2010): 508-525.

[4] Pak, Alexander, and Patrick Paroubek. "Twitter as a corpus for sentiment analysis and opinion mining." LREc. Vol. 10. No. 2010. 2010.

[5] Agarwal, Apoorv, et al. "Sentiment analysis of twitter data." Proceedings of the Workshop on Language in Social Media (LSM 2011). 2011.

[6] Jianqiang, Zhao, Gui Xiaolin, and Zhang Xuejun. "Deep convolution neural networks for Twitter sentiment analysis." IEEE Access 6 (2018): 23253-23260.

[7] Sorensen, David P., and Joseph W. Shepard. "Print-out process and image reproduction sheet therefor." U.S. Patent No. 3,152,904. 13 Oct. 1964.

[8] Dharmarajan, K., and M. A. Dorairangaswamy. "Discovering User Pattern Analysis from Web Log Data using Weblog Expert." Indian Journal of Science and Technology 9 (2016): 42.

[9] Go, Alec, Richa Bhayani, and Lei Huang. "Twitter sentiment classification using distant supervision." CS224N Project Report, Stanford 1.12 (2009).

[10] Tang, Duyu, et al. "Learning sentiment-specific word embedding for twitter sentiment classification." Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Vol. 1. 2014.

[11] Saif, Hassan, Yulan He, and Harith Alani. "Semantic sentiment analysis of twitter." International semantic web conference. Springer, Berlin, Heidelberg, 2012.

[12] Zhou, Xujuan, et al. "Sentiment analysis on tweets for social events." Proceedings of the 2013 IEEE 17th International Conference on Computer Supported Cooperative Work in Design (CSCWD). IEEE, 2013.

[13] Abirami, K., and P. Mayilvaganan. "Fuzzy Clustering with Artificial Bee Colony Algorithm using Web Usage Mining." International Journal of Pure and Applied Mathematics (IJPM)118 (2018): 3619-3626.