

# Hybrid Approach of Multi-Attribute Data for Sentiment Analysis in Online Social Networks

**Dr.Rakesh Kumar Donthi**

Associate Professor, Dept. Of CSE, JNTUH, Geethanjali college of Engineering and Technology,  
Hyderabad, Telangana 501301, India  
[rakesh.cse15@nitp.ac.in](mailto:rakesh.cse15@nitp.ac.in)

**G. Santhoshi**

Assistant Professor, Dept. Of CSE, JNTUH, Geethanjali college of Engineering and Technology,  
Hyderabad, Telangana 501301, India  
[gsanthoshi.cse@gcet.edu.in](mailto:gsanthoshi.cse@gcet.edu.in)

**Maguluri V Lavanya**

Assistant Professor, Dept. Of CSE, JNTUH, Geethanjali college of Engineering and Technology,  
Hyderabad, Telangana 501301, India  
[mlavanya.cse@gcet.edu.in](mailto:mlavanya.cse@gcet.edu.in)

**M.Akhila Reddy**

Assistant Professor, Dept. Of CSE, JNTUH, Geethanjali college of Engineering and Technology,  
Hyderabad, Telangana 501301, India  
[akhilareddy.cse@gcet.edu.in](mailto:akhilareddy.cse@gcet.edu.in)

---

## ABSTRACT

In present days there is an increase in the popularity of MySpace, LinkedIn and other social networks. In these networks there is also increase in the transmission of data among the clients who are available in it. There is an increase in the transmission of data among the clients who are available in it. There is an increase in outsourcing of data and so the data that is moving on social media is also being increased. The data from here is used in different applications and also for research. There are various kinds of existing methods to know sentiment of web based internet Communities (social networks) to mark transmission of data among several users to classify patterns in regard to related attributes to examine the data which is in huge amount. In this paper, we initiate and propose an approach of machine Learning which is Hybrid amalgamation of classification and Balanced Window which is based on Parts Of Speech to operate the data (which is outsourced)of internet Communities such as Facebook and for different blogging services which are trained and sorted the relation basing on aspects of emotions like negative or positive and various relations present in social streams. Our proposed approach performance is immensely near to machine learning and main suitable attributes are recognized arbitrarily and carry out sentiment analysis in various streams of data. The result of the experiments done by us exhibit thorough level of results in classification by comparing the approaches which are existing in the environment of real time.

**Keywords:** Sentiment Analysis; Machine Learning; Hybrid Approach.

---

## 1. Introduction

Out of the most rapidly evolving methods for Internet compatibility as of late is making universe of OSNs (online social networks),for example, LinkedIn, Twitter, Face book, and additionally a large clutch of web blogging (logging) administrations [2]. The confusing estimate of information streaming across interactive institutions has made stabbing for bit which are helpful of grasping in organizations which are informal having a field of great zeal for delayed situations. Due to its extensive size of stream of information, the extraction of information, the extraction of information in interactive institutions has grown to be a familiar field of research, with graded evaluation to be a zone which is outstandingly gripping. The users of interactive institutions can occasionally and one more time be a segment into specific conventions in illumination of planned interims. By identifying these conventions it is designable to exhibit their common evaluation as an agent of a greater community, employing the sub-division in a particular OSN as a check concept investigation analyzes the information shown by individuals in the interior of larger conventions and specified an exemplar, grabs into consideration the guarantee of common temperament or sentiment of that convention regarding particular tips.

In any instance there are distinct strains that are posed by streaming online networking data. The most prominently is the plan of knowledge online networking details can be implied to as “short content information”. The data available is persistently not extremely numerous characters, which shapes much of the subject presentation computations prodigal; as many maxims can’t consistently be decided from such type of information. Another trial is postured by structure of information as such. Late web culture have offered ascend with different colloquialisms as well as tiny structures for instance “Laughing Out Loud”(LOL) and “Talk To You Later”(TTYL) etc. Articulation of any feeling is currently made through emojis. The information which is acquired from online social networking spots are consistently filled with terminologies, hash tags and feelings, in this method building traditional pack of words organizations computations is prodigal.

The excessive accessibility of the user generated subject in social gushes, in any instance, escort a little strains. The broad size of information builds tough to acquire the crucial data in an efficient and potential method. Suggested approaches must be adequately frank to strengthen, yet must control complicated information. A few of these strains are identified with the usual NLP(Natural Language Processing) strategies, for instance, affiliation of supposition highlight [11], invalidation of feeling [2], mockery and incongruity [12][18], and spam conclusion[13]

Interestingly others are recognized with matters for content produced by the client in online for instance, number of dialects, uncommon state of unreliability and ambiguity, wrong spelling, dialects and dirty words[10]. In such unique situation it will be likewise vital to tell the necessity of deciding clients notoriety as well as belief. For specific points lion’s share sentiment (i.e., the shrewdness of the group) may be best arrangement[49] but for others; just specialists suppositions ought to be wellspring of the data to be considered[17]. One more applicable matter is presence of specific sections and types of data that exists within social streams; unequivocal references to the clients, gatherings as well as associations (e.g., @robinwilliams in Twitter), unequivocal structures to allude the ideas (e.g., hash tags of Twitter #funny and # comedian), terms of slang as well as emotion which notice feelings as well as temperaments(lol and D), components which express interests as well as tastes(eg., preferences of Face Book) and the URLs which assets supplement the data which is posted. In view of utilization of the metadata which is logical likewise assumes an important part, extricating and extracting time as well as meta data of geo-area may extremely important to assess investigation on the unique as well as worldwide social stream information. By considering the complexities above which are there in social networks we develop and present a hybrid Machine Learning strategy that is combination of Balanced window which maintains data and classification referencing parts of speech using which outsourced data of social networks like Blogging Services and Face Book is handled. The approach mostly focus on recognizing classification with allowable depiction of various twitter as well as other social networks to categorize various data which is based on various expressions such expressions with regard to neutral, irrelevant and polar data either it will be positive or negative in various frameworks. In addition, the above approach mostly make various area of locality details that are collected through various sources of data like social, railway and more formats of environment which is real time.

## 2. Related Work

The most punctual work on estimation examination on Twitter details accompanied away to the work of Go et al.,(2009) [5]. The process in which they dealt the identification pertaining mockery was to distinguish tweets as +ve or -ve using Support vector machines, Naïve Bayes, Maximum Entropy computations for grouping bigrams as well as unigrams, unigrams along with labels of grammatical form as extractors of highlight. The paper that is proposed ha an arrangement display of emotion identification as well as investigation to arrange as well as further thought to be redundant letters. The computation explained an accuracy around 80%. Be that as it may did not explore mockery recognition, as well as confined setting to +ve and -ve classes. The next stage her of the research will be influencing characterization genuine to time. The ongoing nature of the tweets have inspired various studies on. Among those one being designed by Bifet et al., (2011) [6] that uses streaming of Twitter API to get the tweets continuously. The paper indicates utilization of Massive Online Analysis(MOA); which collects continuous details, uses collection of computations and characterizes tweets to 2 classes, to be specific +ve and -ve tweets. The computation proposed uses a component age channel that uses a plan of weighting as well as performs recognition of change to get in touch at its results. The toil, in any case is again restricted only to 2 class order and will not walk into recognition of mockery. In addition it uses an already defined dataset for preparing. Other look into that is similar to Bifet. A. what will be more, Frank E. (2010) [7] that again uses twitter gushing API as well as already defined to prepare a set for arranging the tweets. The work will assess 3 computations in specific Stochastic Gradient Descent, Hoeffding tree, Multinomial Naïve Bayes to group and show the results of to 82% on best fit computation. Not withstanding, the work which proposed will neglect to do mockery identification, yet uses smiley lexicon to assist in order. Barbosa et al [11] show other approach of supposition order of twitter details in which they named thousand tweets utilizing extremity forecasts from 3 unique sites as well as other thousand Tweets for testing. Those investigated a few qualities of tweets , for example, how they are composed, and additionally meta-data of the phrases which are used to form them. In growth to extremity of terms and

parts of speech of the terms they use sentence structure highlights of tweets, for instance hash tags, retweets, accentuations and marks of outcry. The precision of test outcomes got was large because of (i) the portrayal of messages were made more dynamic as opposed to utilizing crude word portrayal (ii)sensible quality's marks given by sources of information are being consolidated.

Pfitzer, Garas, and Schweitzer [3] ordered the posts of twitter which are serving one in two capacities unmistakably. Creation of data (or) the uncontaminated data appropriation in which a client responds other's unique thought (or) idea (re-tweet). It has been discovered the enthusiastic dissimilarity affects likelihood in a snippet of data which is being reposted; the tweets having more passionate disparity have large likelihood of being reposted. The work done in past in the estimation of examination of twitter has uncovered that there is particular connection of the aggregate state of mind because of capacity to retweet another client's post. In one of the explorations done by Garas, Garcia, Skowron, and Schweitzer [4] the correlated designs utilized crosswise upon online chat rooms would seize up for the examples of the both constant data trade as well as general feeling. Examples which are frequently arise in exchange types of online correlation, quick response looks are more adjusted in understanding about different points. The dialect of the tweets is exceptional because of the 140-character restrict forced upon singular posts, making clients regularly use shorthand documentation and emotions in slant articulation [5]. Paroubek[6] and Pak have done a worthwhile work in short content examination field when grouping the tweets in view of connection amongst's emotions and sections of discourse. The emotions were utilized in deciding general assumption of tweet [7], as far as possible which makes that further probable for it to be a just single conclusion. The rest of it was part in particular sections of discourse using the calculation of Tree Tagger, demonstrating whatever grammatical feature is having the best effect on the post's general slant [6] [8]. Much firmly identified with examination tendered in the present paper is crafted by Aston, Liddle, and Hu [9] which is on assumptions of Twitter grouping information flow. Calculation using perceptron and its chosen rendition in the company of fixed element choice have been utilized to foresee notion present inside an information flow condition in favor of a lot of twitter information. The expectation of them is gushing and accomplished comparable accuracy compared to clump forecast of [2]. The words in the best highlights were revealed.

In the same way we aim to apply the Hybrid Machine Learning criteria to online community feeling category in information flow and want to apply the selection of features online in combination along it so as to consider the modifying information slowly but surely.

### **3. Procedure for Pre Processing of data for Data streams**

To define the data pre processing the basic evaluation method is pre processing. Due to the changing and abnormal character of accent used as a section of tweets, it is probable that pre-processing methods make use of incorporate specific tokens of tweets. It is profoundly probable that many tweets contain few type of linguistic or spelling botches, abbreviations, proverbs, dialects; forced into due to 140 constrain forced by Twitter on tweets. Pre- processing procedure extricates the relevant material from tweets while ignoring the superfluous ones. The strategies attached in this paper are used commonly in data recovery applications particularly in assessment investigation in smaller scale blogging. The gathered information is gone across a development of the pre processor which aid change of strings of message towards the component vector.

A portion in the steps of pre-processing which has been completed are clarified underneath. This is one among the imperative strides within the whole arrangement procedure as nature of highlights/ characteristics which are removed from preparation of dataset using the told pre-processing system specifically influences the execution of the classifiers. Mockery will be identified if the first segment in sentence is sure and second one is negative. The element weight is allocated in between of 0 and 1 considering the above algorithm. A component that says if sentence really passes on an assumption or not and it will do how firmly the feeling is underlined. Another test called subjectivity test is done for the whole sentence and later for both parts to assure weighting made ready suitably. The later component will be considered as capitalization. Tweet that consists a considerable measure of uppercase words needs to be made a solid point or a pass on forceful feeling. This is the way the element will be set in view of quantity of the promoted events in a tweet. The above mentioned labels will be considered and inspiration as well as pessimism highlights will be weighted. The two parts of sentence are rehashed for the same. One more vital element that will be observed in parts of speech labels. The section in the discourse will be labeled in light of an existing Characteristic Language Toolkit NLTK collection and then appropriately weighted. For the sake characterization this is the way in which, out of all twenty two highlights are identified and used. After recognizing and weighing each one of these highlights, one more sketch is given by using each one of these highlights along with their different weights. For testing as well as preparing the calculations as group the portrayal will be encouraged with the help of contribution. This is one among the imperative strides within the whole arrangement procedure as nature of highlights/ characteristics which are removed from preparation of dataset using the told pre-processing system specifically influences the execution of the classifiers.

#### **3.1 Pointers Identification (hash tags and usernames)**

In the Twitter, before any username there is use of a token anL @ which points on different clients. Furthermore, clients label these tweets relating a particular classification in the twitter, utilizing #. Once more, to evade blast of the highlights, it is isolated to a steady image <HASHTAG> and <USER>. The substitution of hash tags as well as usernames lessen element estimate because of huge edge.

### 3.2 Stop Words Expulsion

In Retrieval of information, this is a typical procedure to expel words which are to great degree normal(have high IDF esteem) which will not escalate value of arrangement procedure. The, a, an are the basic words mainly called stop words. Considering such words within a tweet will not give valuable data, and are evicted.

### 3.3 Pressure of the Words

The clients of the Twitter have a habit to be exceptionally casual in their dialect and the greater part of their external words to communicate compelling emotions. For example, the expression “happyyyyyyyyy” bunches a higher articulation degree rather than “happy”. During preparation as well as assessment, of words which contains more than three consequent event of the similar rehashed letter/character, it is lessen to a succession of 3 characters by us. For example, let us consider “cooooooooool” to “cool”. Grouping isn't lessened to the more typical two-character cool keeping in mind the end goal to detach among normal utilization as well as underlined usage of a certain word.

### 3.4 Skewness Evacuation in the Dataset

The time at which we get an imbalanced preparation dataset, assembling it is difficult to undertake the characterizations which are helpful. The unevenness of class shows an affair in use of customary characterization calculations because they endeavor to build representations with the aim of augmenting general order precision. To lighten the affairs connected with irregularity of class numerous procedure have been proposed. For instance boosting and information testing.

## 4. System Implementation

Our procedure for implementation of proposed system contains Preparation of data and Sentiment analysis.

### 4.1 Preparation of Data:

Once the steps of pre-processing are finished, information that has to be arranged and prepared for grouping stage. Tweets are accumulation of the sentences that can't be specifically nourished into classifiers. Consequently five noteworthy advances are carried out in order to set up information for following stage. Different stages are a) Organization of words b) parts of speech(POS) labeling c) Lemmatization and Stemming d) id of the feature e) creation of new representation. Tokenization is the initial step. It is accomplished on the tweets in order to split them to culminate significant segments from a particular sentence. Some of the time tokens can be as far as passages or entire sentences however it is shown as a word in the proposed. Tweet is split into watch words and words that help in the characterization of evacuation of picked as well as stop words. Parts of Speech labeling is performed after tokenization of tweets. Words in the tweet and parts of them of discourse assume part in the characterization. The event which individual is utilizing a considerable measure of modifiers there is a plausibility that he is portraying something with a lot of acclaim, that insights about it being wry. In view of this, the proposed show labels the sections of discourse of a particular word. Lemmatization and stemming is performed after the completion of labeling. The stemming is based on the possibility of words which are having a similar stem and are shut in importance. After this the ID specific weights will be appointed in aspect of significance. The procedure for distinguishing proof of root expression of different words used as part of tweet. For example, the words such as mice will be changed to mouse. This type of transformation illuminates setting of use for word and makes them less demanding for delineation of its importance. After making all the important arrangements for include recognizable proof real highlights are considered which are a)Polarity of Blob b) Subjectivity of blob c) Positive sentiment d)Negative assumption and e) Capitalization. Subject and theme are other couple of highlights that are plain as day. The principal include the extremity of the blob, it suggests the conclusion which is passed on by sentence like entirety. It says whether the sentence is certain, impartial or negative. Such type of highlight is initially connected to entire sentence and next for initial and later parts of sentence. Algorithm for implementation of

classification of data for data preparation

Mockery will be identified if the first segment in sentence is sure and second one is negative. The element weight is allocated in between of 0 and 1 considering the above algorithm. A component that says if sentence really passes on an assumption or not and it will do how firmly the feeling is underlined. Another test called subjectivity test is done for the whole sentence and later for both parts to assure weighting made ready suitably. The later component will be considered as capitalization. Tweet that consists a considerable measure of uppercase words needs to be made a solid point or a pass on forceful feeling. This is the way the element will be set in view of quantity of the promoted events in a tweet. The above mentioned labels will be considered and inspiration as well as pessimism highlights will be weighted. The two parts of sentence are rehashed for the same. One more vital element that will be observed in parts of speech labels. The section in the discourse will be labeled in light of an existing Characteristic Language Toolkit NLTK collection and then appropriately weighted. For the sake characterization this is the way in which, out of all twenty two highlights are identified and used. After recognizing and weighing each one of these highlights, one more sketch is given by using each one of these highlights along with their different weights. For testing as well as preparing the calculations as group the portrayal will be encouraged with the help of contribution.

**Information** : Personally categorized twitter upgrade dataset  
**Result** : Binary category into ironic and non data  
 Initialise slangDictionary, emojiDictionary Import dataset  
**step 1: Information Pre-processing**  
 Replace emoji in twitter upgrade with emojiDictionary value Replace terminology in twitter upgrade with slangsDictionary value  
 Hashtag separation  
**step 2: Information preparation** Word tokenization of pre-processed Tagging the tokenized terms with parts of speech (POS) Stemming and Lemmatization Feature recognition and weighting Update the dataset with the addition of the options and respective weights  
**step 3: data extraction** Classification of ready data using 5 classifiers Choosing the classifier with best accuracy Training the design with best classifier Testing the design with best classifier Real time examining of tweets.

**Algorithm 1. Real time classification model to define data classification.**

**4.2 Sentiment Analysis:** In social networks depiction of real time of sentiment analysis the balanced window(BW) is used for calculation of sentiment analysis over the whole data depiction of the network. BW enquires an advancement variable  $\alpha$  as well as downgrade variable  $\beta$ . This is isolated from the Balanced window on using bigger edges as well as a slight change in refresh rules for weight frameworks. BW is illustrated in particular in the figure Balanced Window.

We have initially prepared BW classifier on the first hundred occasions. For every occurrence we initially group it by using MBW, and if sequence was right we go to the next case and refreshes the correct. If it was mistaken that we have refreshed the off base check and try to refresh the weight of framework.

1. Initialize  $I = 0$ , reverse  $c = 0$ , and designs  $u_0$  and  $v_0$   
 2. For  $t = 1, 2, \dots, T$ :  
 a. Get new example  $x_t$  and add prejudice.  
 b. Stabilize  $x_t$  to 1.  
 c. Determine ranking  $= (x_t, u_i) - (x_t, v_i) - \theta t$   
 d. Get real category  $y_t$   
 e. If forecast was wrong:  
 i. Upgrade Models. For each function  $j$  where  $x_t > 0$ :

**Algorithm 2. Balanced window representation for sentiment data analysis.**

## 5. Experimental Results

The approach that is proposed needs 4 different client characterized variables like (great component choice, downgrade, advancement, include choice), which might be prompting reduction in the exactness while picked randomly. This is the way we perform broad checks on these variables to reveal feasible important reaches of each. It was solved that more accurate our forecast is progressed to becoming, the bigger is gram measure used, because the big gram sizes reveal much a word other than littler gram magnitudes. Further, while using highlight establishment as well as choice of the great element. It can be observed that accuracy will be most notable with many as well as great highlights used as seen in Table 1

Table 1. Accuracy of sentiment prediction with different data grams.

	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
0.1	58.3	61.6	64.2	65.4	66.6	67.4	66.9	68.8	69.8	71.8
0.2	58.9	61.7	64.4	65.4	67.1	66.7	68.4	68.6	69.8	72.6
0.3	60.2	62.9	64.4	66.5	67.5	67.7	68.9	68.4	69.4	71.9
0.4	59.6	62.7	64	66.1	66.8	67.6	68.8	69.8	70.3	73.3
0.5	60.9	62.5	65.1	66.1	66.9	67.9	68.4	70.7	69.3	71.7
0.6	61.2	63.6	65	64.5	67	68.1	69.9	68	68.8	72.6
0.7	61.4	63.2	64.3	65.3	68	69.1	67.3	69	69.9	72.2
0.8	63.4	65.2	64.5	65.5	67.7	66.3	68.4	69.2	69.3	72.6
0.9	64.8	66.5	66.7	65.4	68.2	65.6	67.8	68.3	69.5	72
1	67.4	67.4	67.4	67.4	67.4	67.4	67.4	67.4	67.4	67.4

We have separated our datasets in consecutive portions of proportions hundred to uncover significance of the highlights after some time. We at that point ran BW on each section of 100 as a flow.

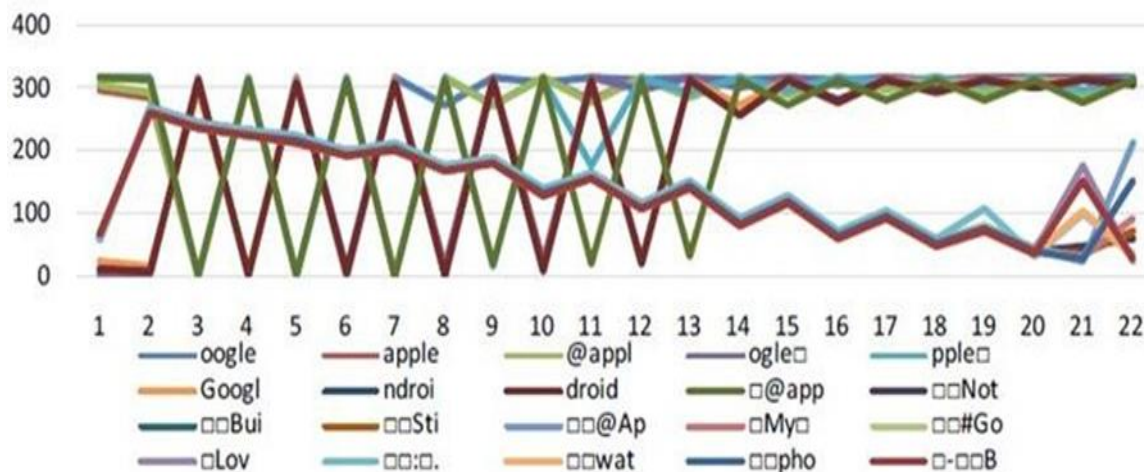


Figure 1. Performance evaluation with respect to sentiment analysis in real time scenario.

The above figure delineates element establishment of main twenty includes over twenty two timestamps of Sanders 5 gram portrayal. The above figure shows, best highlights waves of significance at outset, yet as information will be ceaselessly encouraged in, few best highlights start to balance out as well as hold the higher position. The highlights which are subsets of words "Google", "Apple" as well as "Android" are some highlights which hold higher significance afterwards. For every particular gram size in Sanders, high precision expectation will fall inside scope of the advancement as well as downgrade estimations of 0.9 as well as 1.1. A top precision will be seen while  $\alpha$  and  $\beta$  values increment concurrently from about  $\{1,1\}$  ahead. When the choice of dynamic element is consolidated in the MBW, we

have achieved exactness with 73.3% while [9] achieved a precision with 77% using five grams with a choice of manual component. As far as an information stream, this is vital to do determination of dynamic element because of the changing importance of the highlights in the company of new approaching information

## 6. Conclusion

We present as well as propose a Hybrid Machine Learning Approach so as to evaluate the Sentiment Analysis on data of real time of various applications. Due to brisk rise of social networks in the environment of real time in this paper. We investigated the matter of idea examination and speculation grouping o social corporation little blog details which as spoke regarding is totally not same like other order of estimation matter on sorted out as well as sure messages. We have broken down the pre-paring of unprocessed messages of twitter (tweets) in particular and gave rules to make reasonable groups for making, in light of research writing. Proposed approach gives efficiently exact results to the both immature/impartial/polar as well as negative/positive preparing groups, while essential Naïve Bayesian Classifier has neglected to do 2 cases. Some more challenging difficulties in organic processing of language can be used as next additions of the above study, like recognition of sarcasm, evaluation managing, context switches. The world wide expression category and terms in foreign can be also researched in added further details in future.

## References

- [1] Anukarsh G Prasad; Sanjana S, Skanda M Bhat, B S Harish, “Sentiment Analysis for Sarcasm Detection on Streaming Short Text Data”, 2017 2nd International Conference on Knowledge Engineering and Applications.
- [2] Nathan Aston, Timothy Munson, Jacob Liddle, Garrett Hartshaw, Dane Livingston, Wei Hu\*, “Sentiment Analysis on the Social Networks Using Stream Algorithms”, Journal of Data Analysis and Information Processing, 2014, 2, 60-66.
- [3] BalakrishnanGokulakrishnan \* 1, Pavalanathan Priyanthan \*2, ThiruchittampalamRagavan, “Opinion Mining and Sentiment Analysis on a Twitter Data Stream”, The International Conference on Advances in ICT for Emerging Regions - ICTer 2012 : 182-188.
- [4] Hassan Saif, F. Javier Ortega, Miriam Fern´andez, Iv´an Cantador, “Sentiment Analysis in Social Streams”, In: Proceedings of the 1st European Conference on Social Media (ECSM’14), pp. 174–182
- [5] Go, A., Huang, L., Bhayani, R.: Twitter sentiment classification using distant supervision. In: CS224N Project Report, Stanford (2009)
- [6] Bifet, A., Holmes, G., Pfahringer, ., Gavald’a., R. “Detecting Sentiment Change in Twitter Streaming Data”, Workshop and Conference Proceedings 17 (2011) 5–11, 2nd Workshop on Applications of PatternAnalysis.
- [7] Bifet A., and Frank E. Sentiment knowledge discovery in twitter streaming data. In Discovery Science, pages 1– 15, 2010.
- [8] Buscaldi, D., Rosso, P, Reyes, A., “From humor recognition to irony detection: The figurative language of social media”, Data & Knowledge Engineering, April 2012