

Sentiment Based Analysis on Twitter Dataset: A Classification Technique

Sonu Patidar

M.tech Student, AITR, RGPV, Bhopal, Madhya Pradesh, India

Abstract -Social networks have revolutionized the way in which people communicate. Information available from social networks is beneficial for analysis of user opinion, looking at the response to policy change or the enjoyment of an ongoing event. Manually sifting through this data is tedious and potentially expensive. Sentiment analysis is a relatively new area, which deals with extracting user opinion automatically. Emotion sensing or sentiment analysis is a complicated task of machine learning technology. That belongs to the Natural language processing branch of artificial intelligence. It is a broad domain of learning and analysis. Among the text classification and their emotional orientation discovery is also part of this domain. In this paper, we proposed decision tree algorithm based text classification model for performing sentiment on twitter based data. Additionally the comparative performance is also measured with the traditional ID3 algorithm and similar variant of the improved ID3 classification algorithm. In order to compare the performance of the algorithms the accuracy, error rate, memory consumption and time consumption is taken as stand parameters.

Keywords: ID3, Sentiment, Tweeter, Social Media, Data Mining, Text Mining, NLP.

I. INTRODUCTION

Twitter [1] is a social networking application which allows people to micro-blog about a broad range of topics. Micro-blogging is defined as "a form of blogging that lets you write brief text updates about your life on the go and send them to friends and interested observers via text messaging, instant messaging (IM), email or the web. Twitter helps users to connect with other Twitter users around the globe. Successful micro-blogging services such as Twitter have become an integral part of the daily life of millions of users. In addition to communicating with friends, family or acquaintances, micro-blogging services are used as recommendation services, real-time news sources and content sharing venues.

These tweets tend to spread to a large number of users in very little time. Users on Twitter not only tweet about their personal issues or nearby events, but also about more general topics or news [2]. Due to the large amounts and diversity of real-time information contained on the site, Twitter lists a freshly updated set of trending topics. Trending topics comprise the top terms being discussed currently on Twitter. This list of top terms, which is updated in real-time, provides a reflection of the current

main interests of the community, i.e., the most-discussed conversations right at the moment.

A. Characteristics of Tweets

Twitter messages have many unique attributes, which differentiates our research from previous research [10]:

Length The maximum length of a Twitter message is 140 characters. From our training set, we calculate that the average length of a tweet is 14 words or 78 characters. This is very different from the previous sentiment classification research that focused on classifying longer bodies of work, such as movie reviews.

Data availability another difference is the magnitude of data available. With the Twitter API, it is very easy to collect millions of tweets for training. In past research, tests only consisted of thousands of training items.

Language model Twitter users post messages from many different media, including their cell phones. The frequency of misspellings and slang in tweets is much higher than in other domains.

Domain Twitter users post short messages about a variety of topics unlike other sites which are tailored to a specific topic. This differs from a large percentage of past research, which focused on specific domains such as movie reviews.

The use of digital text is increasing as the social media increases their effect in daily life. A number of research groups and individual researchers are working to find the patterns on these data. In this study the social media text analysis and sentiment analysis techniques are investigated and a new classification technique is proposed for enhancing the performance of text classification. The given chapter provides an overview of the proposed work and involved investigation.

B. Sentiment Analysis

Sentiment analysis can be defined as a process that automates mining of attitudes, opinions, views and emotions from text, speech, tweets and database sources through Natural Language Processing (NLP). Sentiment analysis involves classifying opinions in text into categories like "positive" or "negative" or "neutral". It's also referred as subjectivity analysis, opinion mining, and appraisal extraction. The words opinion, sentiment, view

and belief are used interchangeably but there are differences between them.

- ✓ *Opinion*: A conclusion open to dispute (because different experts have different opinions)
- ✓ *View*: subjective opinion
- ✓ *Belief*: deliberate acceptance and intellectual assent
- ✓ *Sentiment*: opinion representing one's feelings

Sentiment Analysis is a term that includes many tasks such as sentiment extraction, sentiment classification, and subjectivity classification, summarization of opinions or opinion spam detection, among others. It aims to analyze people's sentiments, attitudes, opinions emotions, etc. towards elements such as, products, individuals, topics, organizations, and services [3].

C. Feature Extraction

The preprocessed dataset has many distinctive properties. In the feature extraction method, we extract the aspects from the processed dataset. Later this aspect are used to compute the positive and negative polarity in a sentence which is useful for determining the opinion of the individuals using models like unigram, bigram [4]. Machine learning techniques require representing the key features of text or documents for processing. These key features are considered as feature vectors which are used for the classification task. Panget al. [5] showed better results by using presence instead of frequencies.

- ✓ *Parts of Speech Tags*: Parts of speech like adjectives, adverbs and some groups of verbs and nouns are good indicators of subjectivity and sentiment. We can generate syntactic dependency patterns by parsing or dependency trees.
- ✓ *Opinion Words and Phrases*: Apart from specific words, some phrases and idioms which convey sentiments can be used as features. E.g. cost someone an arm and leg.
- ✓ *Position of Terms*: The position of a term with in a text can effect on how much the term makes difference in overall sentiment of the text.
- ✓ *Negation*: Negation is an important but difficult feature to interpret. The presence of a negation usually changes the polarity of the opinion. e.g., I am not happy
- ✓ *Syntax*: Syntactic patterns like collocations are used as features to learn subjectivity patterns by many of the researchers.

II. LITERATURE SURVEY

Agarwal et al. [7] present an analysis in which twitter is different to other forms of raw data which are used for sentiment analysis as sentiments are conveyed in one or two sentence blurbs rather than paragraphs. Twitter is much more informal and less consistent in terms of

language. Users cover a wide array of topics which interest them and use many symbols such as emoticons to express their views on many aspects of their life When using human generated status updates, sentiment are not always obvious; many tweets are ambiguous and can use humors to maximize the opinion to other human readers but deflect the opinion to a machine learning algorithm.

ApoorvAgarwal et al [8] examine sentiment analysis on Twitter data. The contributions of this paper are: (1) first introduce POS-specific prior polarity features. (2) And explore the use of a tree kernel to obviate the need for tedious feature engineering. The new features (in conjunction with previously proposed features) and the tree kernel perform approximately at the same level, both outperforming the state-of-the-art baseline.

The influence of micro blog on information transmission is becoming more and more obvious. By characterizing the behavior of following and being followed as out-degree and in-degree respectively, a micro blog social network was built in this paper. It was found to have short diameter of connected graph, short average path length and high average clustering coefficient. The distributions of out-degree, in-degree and total number of micro blogs posted present power-law characters. The exponent of total number distribution of micro blogs is negatively correlated with the degree of each user. With the increase of degree, the exponent decreases much slower. Based on empirical analysis, Qiang Yan et al [9]proposed a social network based human dynamics model in this paper, and pointed out that inducing drive and spontaneous drive lead to the behavior of posting micro blogs. The simulation results of model match well with practical situation.

Another consideration when using a dataset generated from Twitter is that a considerably large amount of tweets which convey no sentiment such as linking to a news article, which can lead to difficulties in data gathering, training and testing. Movassate et al. [10] provides Sentiment analysis of tracking opinions and attitudes on the web and determines if they are positively or negatively received by the public.

III. PROPOSED WORK

A. Methodology

The proposed system for the sentiment text analysis and their accurate evaluation a new system is prepared using the traditionally available techniques. The organization of traditional methodologies for obtaining sentiment based text analysis is given using figure 1.

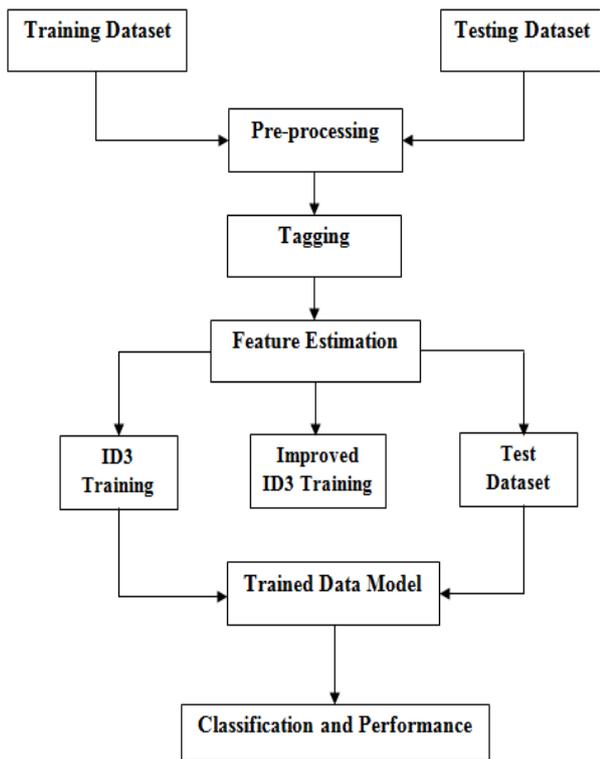


Figure1: Proposed Models

Pre-processing: Both training and test data is pre-processed in this phase, the pre-processing of data involves the removal of punctuations and removal of frequently occurred words sometimes termed as stop word. These stop words is frequently used in each sentences such as is, am, are, this, that, his, her and others.

Tagging: After pre-processing of data it is required to involve the features on data. Therefore the user input tags are involved with the text such as:

Ram is a good boy.

Can be converted into:

Noun adjective noun

Features estimation: After the tagging the original data is converted into a new encoded format. Therefore the tagged data and the associated tag are stored on a relational data base which contains the encoded attributes and their class labels. The example of the text to feature extraction is given using below given table:

Table1: Feature Data

No	Pr	Ve	ad	A	Pr	Co	Go	ba	Class
un	o-	rb	v	dj	e	nj	od	d	es
2	1	1	0	0	0	0	1	0	1

ID3 Training: Given table data used to learn with traditional ID3. Basically here a provision is made to select

algorithm for training. If user select traditional ID3 model the system make training using the traditional ID3.

Improved ID3 Training: If user selects the improved ID3 then the data is accepted through the improved ID3 and developed the decision tree data model form the training samples.

Trained data model: Trained model is a finalized decision tree that makes use the input data and converted into a tree structure. For classification of data test data instance is invoked through the decision tree model and their class labels are predicted.

Test data: That is a part of training sample which is used to perform testing of trained data model using the cross validation technique. The cross validation results the accurate amount of data that is correctly recognized using the decision tree.

Classification and performance: Finally the performance of the entire system is computed in terms of accuracy, error rate, time consumption, and the memory consumption during training and testing of data.

Proposed ID3 algorithm

Input: the training Data D

Output: Rules Set R

Process:

1. $T_r = readTrainingSet(D)$
2. $T_{model} = ID3.Train(T_r)$
3. $E_{rules} = extractRules(T_{model})$
4. *for* ($i = 0; i \leq D.attributes; i++$)
 - a. $F_A = findUniqueAttributeValues(D)$
 - b. *for* ($j = 0; j \leq F_A.size; j++$)
 - i. $P_j = Bays.Probability(D, F_A[j])$
 - c. *End for*
5. *End for*
6. $R = prunRules(E_{rules}, P_j)$
7. *return R*

IV. RESULT ANALYSIS

Accuracy

In a classification technique the accuracy is measurement of accurately classified patterns over the total input patterns produced for classification. Therefore that can be a measurement of successful training of the classification algorithm. The accuracy of the ID3 can be evaluated using the following formula:

$$Accuracy = \frac{Total\ correctly\ classified\ patterns}{Total\ input\ patterns\ to\ Classify} \times 100$$

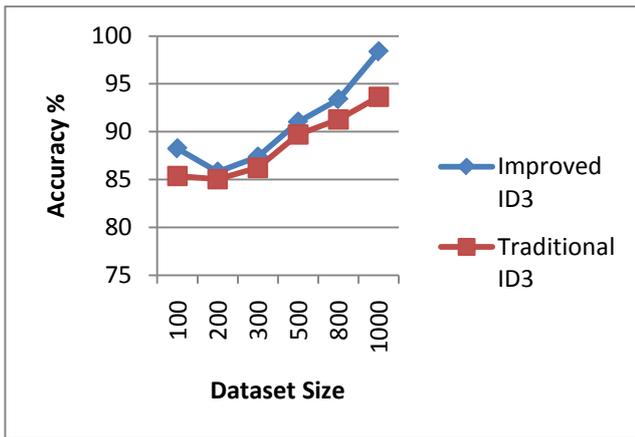


Figure 2: Accuracy

The accuracy of the implemented Improved ID3 is represented using figure 2. The figure contains the accuracy of the implemented algorithms. The X axis of the diagram contains the amount of data during the training and testing and Y axis contains the obtained performance in terms of accuracy percentage. To demonstrate the performance of both the techniques the blue line is used for proposed model and red line shows the performance of traditional ID3. According to the obtained results the performance of the proposed ID3 technique provides more accurate results. Additionally the accuracy of the learning model is increases as the amount of instances for the learning of algorithm is increases.

Error Rate

The amount of data misclassified samples during classification of algorithms is known as error rate of the system. That can also be computed using the following formula.

$$Error\ Rate\ \% = \frac{Total\ Misclassified\ Patterns}{Total\ Input\ Patterns} \times 100$$

Or

$$Error\ Rate\ \% = 100 - Accuracy$$

The figure 3 and shows the comparative error rate of implemented ID3. In order to show the performance of the system the X axis contains the amount of data used for training and the Y axis shows the performance in terms of error rate percentage. The error rate of the traditional ID3 is given using the red line and the performance of the Improved ID3 is given using the blue line. The performance of the proposed classification is effective and efficient during different execution and reducing with the amount of data increases. Thus the presented classifier is more efficient and accurate than the traditional approaches of text classification.

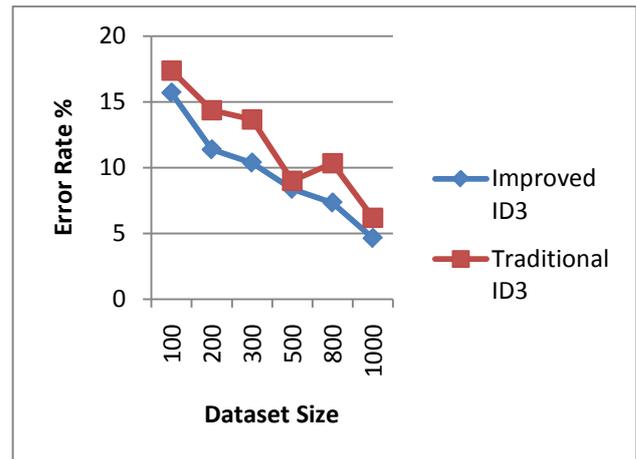


Figure 3: Error Rate

Memory Usage

Memory consumption of the system also termed as the space complexity in terms of algorithm performance. That can be calculated using the following formula:

The amount of memory consumption depends on the amount of data reside in the main memory, therefore that affect the computational cost of an algorithm execution. The performance of the implemented classifier for sentiment classification is given using figure 5.3. For reporting the performance the X axis of figure contains the amount of data required to execute using the algorithms and the Y axis shows the respective memory consumption during execution in terms of kilobytes (KB). According to the obtained results the performance of algorithm demonstrates similar behavior with increasing size of data, but the amount of memory consumption is increases with the amount of data.

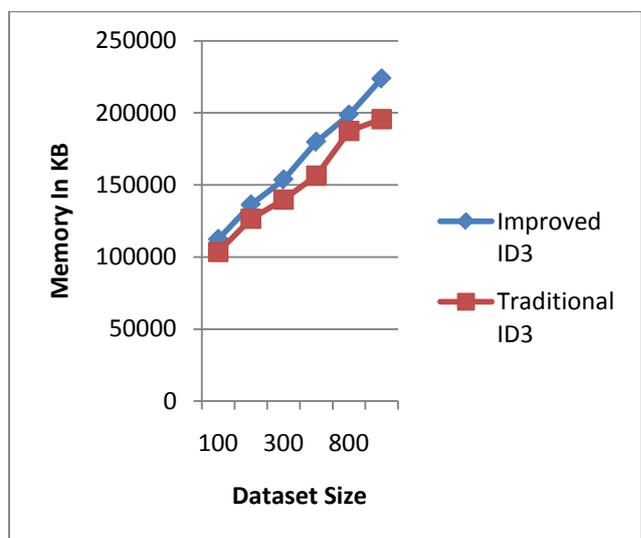


Figure 4: Memory Consumption

Time Utilization

The amount of time required to classify the entire test data is known as the time consumption. That can be computed using the following formula:

$$\text{Time Consumed} = \text{End Time} - \text{Start Time}$$

The time consumption of the proposed algorithm is given using figure 4. In this diagram the X axis contains the size of dataset and the Y axis contains time consumed in terms of milliseconds. According to the comparative results analysis the performance of the proposed technique shows the high time consumption. But the amount of time is increases in similar manner as the amount of data for analysis is increases.

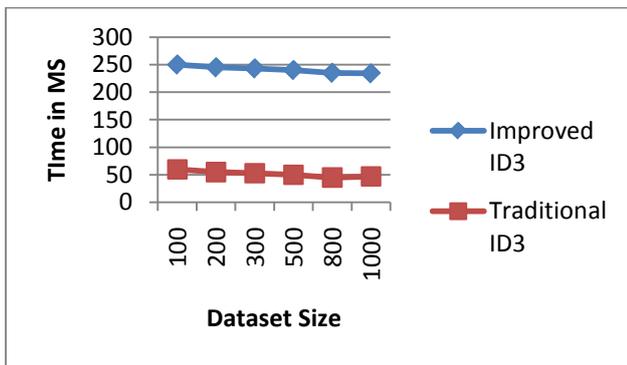


Figure 5: Time Consumption

V. CONCLUSION

Data mining offers the supervised and unsupervised learning concept to analyses the data and classifies or categorize in a predefined groups of data. The algorithms enable us to use the computer based algorithms to analyze the data automatically without any human efforts. The proposed work analyzes the social network based text for their sentiments and orientation based text classification. Therefore the proposed work involves the pre-processing, tagging, learning and the classification of newly arrived patterns. The performance of the system is estimated for finding the system accuracy and error rate for the sentiment data.

VI. FUTURE SCOPE

The proposed work is adoptable and efficient for classifying the patterns of the text data for analyzing emotions hidden in text. Therefore the proposed technique involves a number of techniques for improving the classification rate. In near future that is promising approach for providing the efficient and accurate classification. But need to adopt more literature for improve the technique for big data environment where actually streamed data appeared on the storage. Additionally need to improve the real time data opinion of the user. That approach is also extendable with the implementation of different other kinds of applications

some of the essential applications such as stock market forecasting, multi-labeled classes based classification and other.

REFERENCES

- [1] <http://www.twitter.com>
- [2] E. Mishaud. Twitter: Expressions of the whole self, Master's thesis, Department of Media and Communications, University of London, 2007.
- [3] Vishal A. Kharde and S.S. Sonawane, "Sentiment Analysis of Twitter Data: A Survey of Techniques", International Journal of Computer Applications (IJCA) Volume 139 – No.11, April 2016
- [4] Socher, Richard, et al. "Recursive deep models for semantic compositionality over a sentiment Treebank" Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP) 2013
- [5] Pablo Gamallo, Marcos Garcia, "Citius: A Naive-Bayes Strategy for Sentiment Analysis on English Tweets", 8th International Workshop on Semantic Evaluation (SemEval 2014), Dublin, Ireland, Aug 23-24 2014, pp 171-175.
- [6] Go, Alec, Richa Bhayani, and Lei Huang. "Twitter sentiment classification using distant supervision." CS224N Project Report, Stanford 1 (2009): 12.
- [7] Agarwal, Apoorv, et al. "Sentiment analysis of twitter data," Proceedings of the workshop on languages in social media, Association for Computational Linguistics, 2011.
- [8] Adam Marcus, Michael S. Bernstein, Osama Badar, David R. Karger, Samuel Madden, Robert C. Miller, "TwitInfo: Aggregating and Visualizing Microblogs for Event Exploration", CHI 2011, May 7–12, 2011, Vancouver, BC, Canada. ACM
- [9] Junghoon Chae, Dennis Thom, Yun Jang, SungYe Kim, Thomas Ertl, David S. Ebert, "Public behavior response analysis in disaster events utilizing visual analytics of microblog data", & 2013.
- [10] Parikh, Ravi, and Matin Movassate. "Sentiment analysis of user-generated twitter updates using various classification techniques." CS224N Final Report (2009): 1-18.
- [11] Wilson, Theresa, Janyce Wiebe, and Paul Hoffmann, "Recognizing contextual polarity in phrase-level sentiment analysis", Proceedings of the conference on human language technology and empirical methods in natural language processing, Association for Computational Linguistics, 2005
- [12] Vishal Gupta, Gurpreet S. Lehal, "A Survey of Text Mining Techniques and Applications", Journal of Emerging Technologies in Web Intelligence, Vol. 1, No. 1, August 2009
- [13] P. Bhargavi, B. Jyothi, S. Jyothi, K. Sekar, "Knowledge Extraction Using Rule Based Decision Tree Approach", IJCSNS International Journal of Computer Science and Network Security, VOL.8 No.7, July 2000
- [14] Umajancy. S, Dr. Antony Selvadoss Thanamani, "An Analysis on Text Mining –Text Retrieval and Text Extraction", International Journal of Advanced Research in Computer and Communication Engineering, Vol. 2, Issue 8, August 2013
- [15] H. Karanikas, C. Tjortjis, and B. Theodoulidis, "An approach to text mining using information extraction," in Proceedings of Workshop of Knowledge Management:

Theory and Applications in Principles of Data Mining and Knowledge Discovery 4th European Conference, 2000.

- [16] R. Hale, "Text mining: Getting more value from literature resources," *Drug Discovery Today*, vol. 10, no. 6, pp. 377–379, 2005.
- [17] S. Chakrabarti, *Mining the Web: Analysis of Hypertext and Semi Structured Data*. San Francisco, CA: Morgan Kaufman, 2000.