

Polarity-based Classification of Subjective Micro-posts using Machine Learning Technique

Arshad Arain¹, Rajesh kumar², Nudra Siddiquie³, Komal Naz⁴, Sabeen gul⁵, Shahbaz Wahab⁶

^{1,2,5}Department of Computer Engineering (IICT) Mehran University of Engineering Sciences and Technology

^{4,6}Department of Software Engineering/IT (IICT) Mehran University of Engineering Sciences and Technology

³Department of Computer Engineering Mohamm Ali Jinnah University Karachi

E-MAIL: swe.arshad@gmail.com, rajesh93_kh@live.com, nudrasiddiqui@gmail.com, komalsoonro525@gmail.com
Sabeenmemon22@gmail.com, shahbazwahab143@gmail.com,

Abstract— This paper covers the two methodologies for notion investigation: I) vocabulary based strategy; ii) machine learning technique. We depict a few methods to execute these methodologies and talk about how they can be embraced for supposition order of Twitter messages. We exhibit a relative investigation of various dictionary mixes and demonstrate that upgrading assumption vocabularies with emojis, contractions and internet based life slang articulations expands the exactness of vocabulary based order for Twitter. We talk about the significance of highlight age and highlight determination forms for machine learning assumption characterization. To measure the execution of the fundamental feeling examination techniques over Twitter we run these calculations on a benchmark Twitter dataset from the SemEval-2013 rivalry, undertaking 2-B. The outcomes demonstrate that machine learning strategy in view of SVM and Naive Bayes classifiers beats the vocabulary technique. We display another outfit technique that uses a dictionary based estimation score as info include for the machine learning approach. The joined technique demonstrated to create more exact arrangements.

Index Terms— social media, Twitter, natural language processing, lexicon, emoticons.

I. INTRODUCTION

Conclusion examination is a region of research that explores individuals' sentiments towards various matters: items, occasions, associations (Bing, 2012). The part of conclusion examination has been developing essentially with the quick spread of informal organizations, microblogging applications and discussions. Today, relatively every site page has an area for the clients to leave their remarks about items or administrations, and offer them with companions on Facebook, Twitter or Pinterest - something that was unrealistic only a couple of years prior. Mining this volume of feelings gives data to understanding aggregate human conduct and it is

of profitable business intrigue. For example, an expanding measure of confirmation brings up that by breaking down supposition of web based life content it may be conceivable to anticipate the span of the business sectors (Bollen et al., 2010) or joblessness rates after some time (Antenucci et al., 2014). A standout amongst the most well known microblogging stages is Twitter. It has been developing consistently throughout the previous quite a while and has turned into a gathering point for a different scope of individuals: understudies, experts, big names, organizations and government officials. This notoriety of Twitter brings about the tremendous measure of data being gone through the benefit, covering an extensive variety of points from individuals prosperity to the conclusions about the brands, items, legislators and get-togethers. In this settings Twitter turns into a intense instrument for expectations. For instance, (Asur and Huberman, 2010) could anticipate from Twitter investigation the measure of ticket deals at the opening end of the week for films with 97.3% exactness, higher than the one accomplished by the Hollywood Stock Trade, a known expectation device for the motion pictures.

Understanding user behavior in virtual networked environments (e.g., social media) and its implications. Enhancing knowledge discovery and content analysis through information visualization

II. MACHINE LEARNING APPROACH

- (i) Relay on statistical algorithm to learn automatically the language.
- (ii) Accuracy of results depend on volume of training data.
- (iii) Training set is manually prepared by labeling the texts into pre-defined classes, which serves as an input for the classifier to learn.

II. PROBLEM STATEMENT

- Identifying polarity of micro-posts is an open challenge, manually it is difficult, therefore an

automated solution is needed.

- What if you want to purchase Smartphone
- As it's very expensive
- First you try to get opinion from your friends , colleagues , neighbors etc...
- Still you can't get enough opinion
- May be because of a very fewer people have smart phones experience
- Or can be other reason too
- Search on internet to go through 1000s of web pages & reviews
- What is other useful option
- Sentiment Analysis

III. OBJECTIVES

- To investigate research and work done in sentiment analysis for medical domine
- To build the required dataset by extracting subjective twitter feeds regarding smart phones using Twitter API and weeka software package
- To classify the polarity of tweets in positive and negitive using one of the supervised classification technique to find out over all sentiment orientation of smart phones
- Crawling micro-posts (i.e., Tweets)
- Exploration of machine learning algorithms (i.e., Neural Networks, K-NN, and Naive Bayes)
- Building Classification Model
- Evaluation of the proposed Model

IV. RELATED WORK

In [2] paper, apply an opinion mining approach to summarize the unstructured and ungrammatical users' reviews, based on Support Vector Machine (SVM).

Two levels of classification is applied:

- 1) Features classification
- 2) Polarity classification for every feature class

an opinion mining approach was proposed to mine unstructured and ungrammatical customers' reviews. It was based on splitting the product's reviews into a collection of sentences.in this paper author used 'Linear as well Classification' technique.

In [3] paper The results shows that SVM gives better performance. Author have done twitter mining to predict reputation of Different Mobile Companies. author used 'Linear & probabilistic' Technique.

In[4] paper author did comparison of lexicon-based approaches for sentiment classification of micro-blog posts lexicon-based approaches for sentiment analysis methodology used to evaluate results and performances .

In[5] paper The hybrid Particle Swarm Optimization (PSO) method was used to improve the opinion analysis of best

movie of the year. This study has shown that PSO affect the accuracy of SVM after the hybridization of SVM-PSO. The best accuracy level that gives in this study is 77% and has been achieved by SVM-PSO after data cleansing. Linear & Hybrid Approach used to achieve the desired goal.

In [5] paper using the corpus, they build a sentiment classifier, that is able to determine positive, negative and neutral sentiments for a document. From the observations, conclude that use of syntactic structures to describe emotions or state facts. Some POS-tags may be strong indicators of emotional text. Lexicon-based methodology used to evaluate ths results.

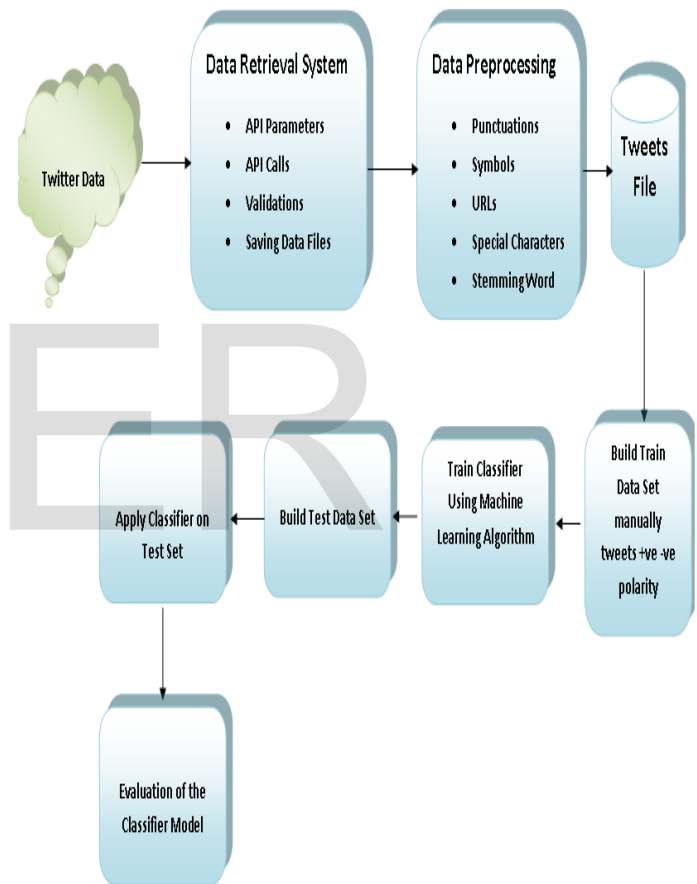


Fig.1 Block Diagram of research methodology

V. METHODOLOGY

In this paper, Classification of "subjective micro-post using Supervised

Machine Learning Technique" to automatically provide concise summary result with best possible methods.

-KNN(K Nearest Neighbors) : Non parametric Function

-SVM(Support Vector Machine): Linear Family

-NN(Neural Network): Non Linear Family

-Naive Bayes: Probabilistic Function

Fig.3 Preprocessing and xlsx file generation

(A) DATASET PREPARATION

STEP 1: TWEET EXTRACTION

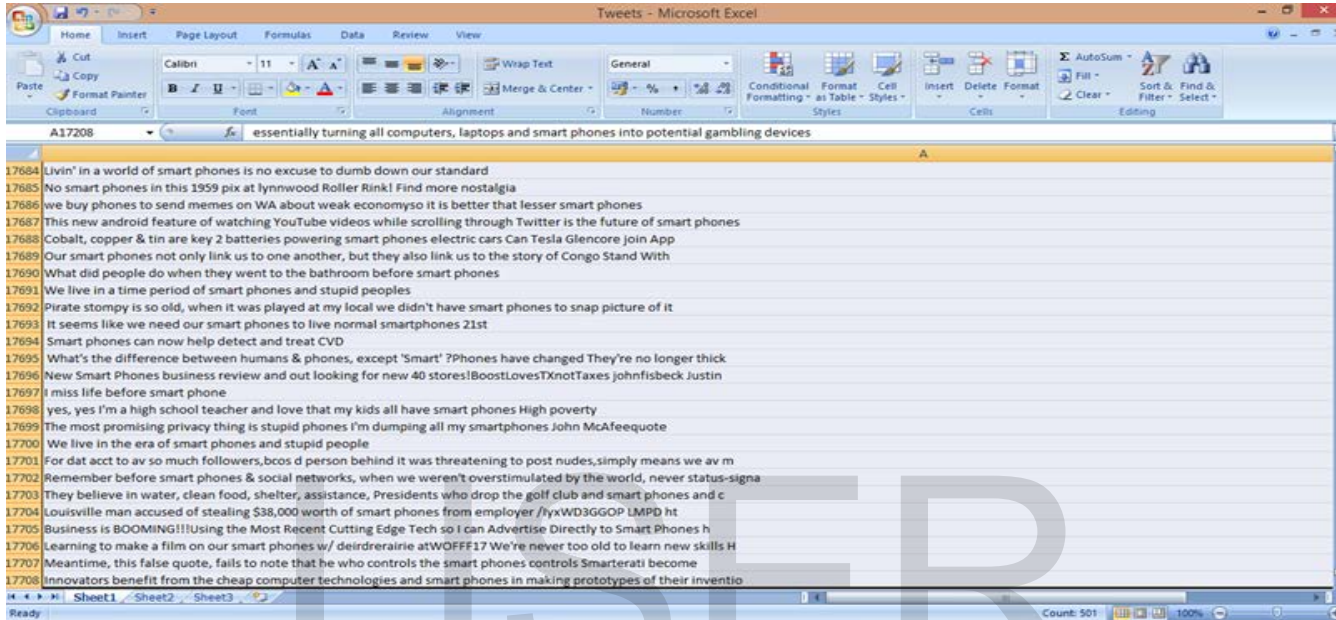
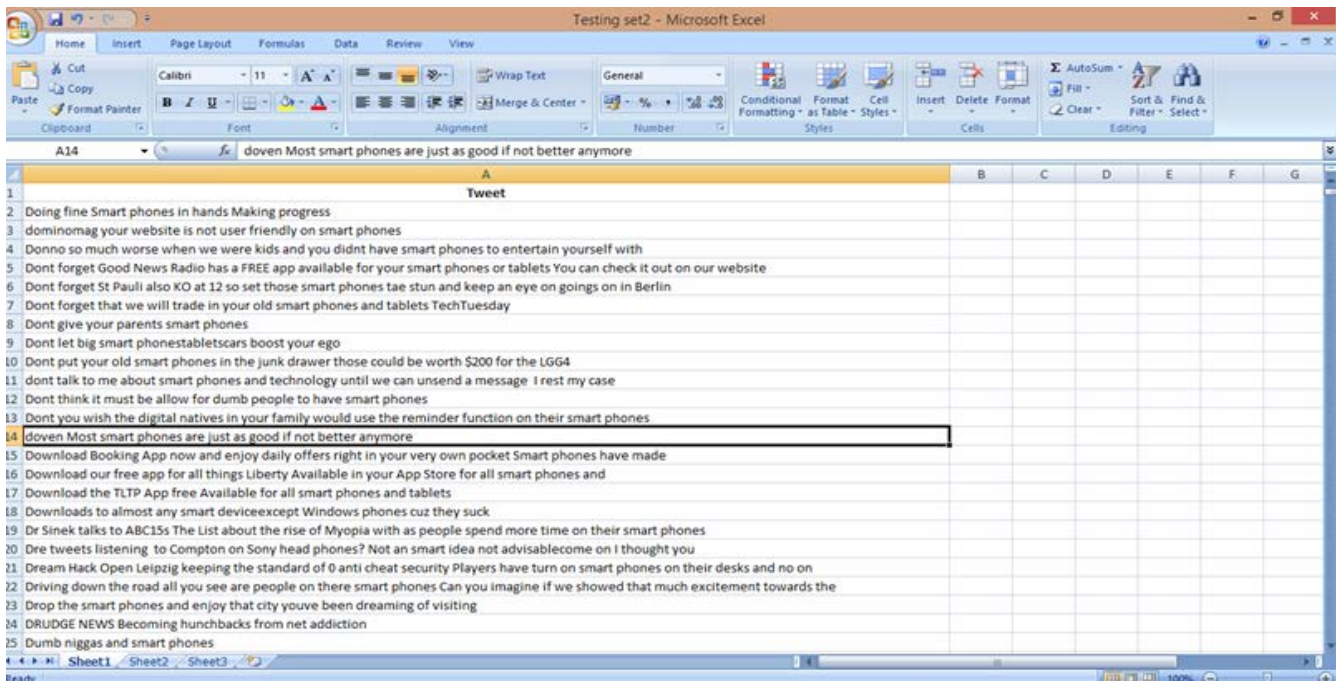


Fig.2 no. of tweets extracted using different keywords

Fig 2 showing the step 1 ,tweet extraction by using different kewywords to evaluate result.

STEP 2: TWEET EXTRACTION



Step 3 .Traning Phase

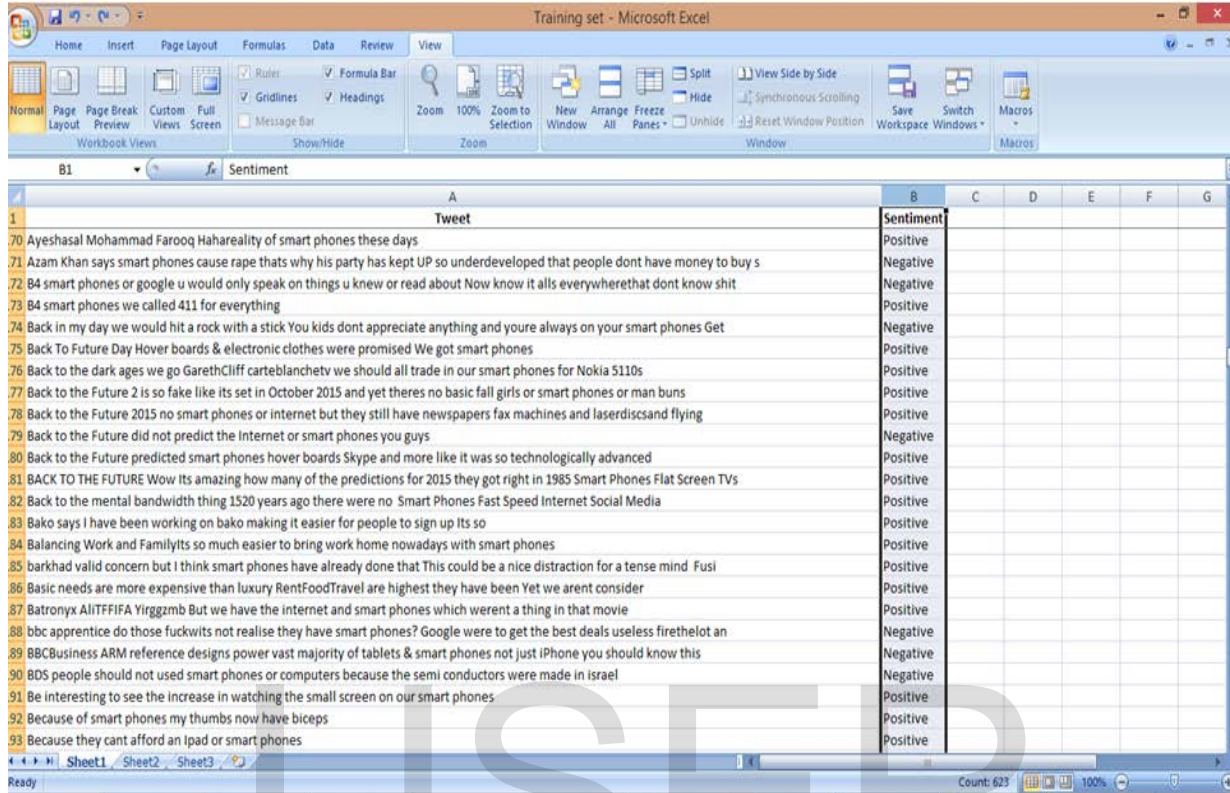
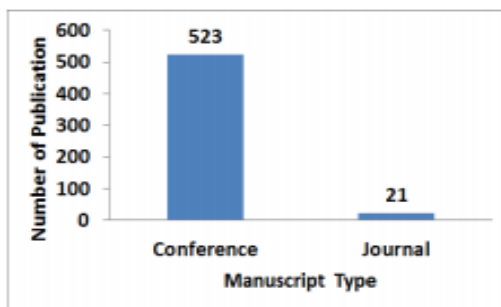


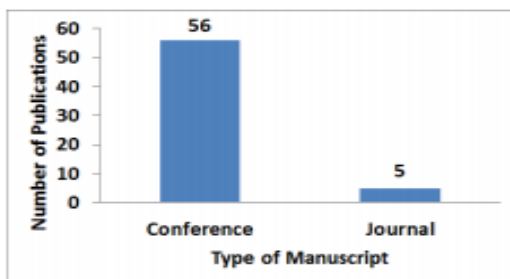
Fig.4 Traning Phase



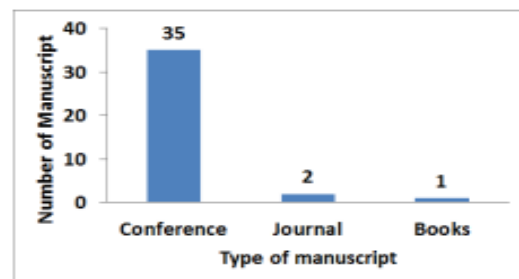
(a) Twitter data



(b) Semantics



(C) Ontology



(d) semantic and ontology

Fig.5 Research trends towards problems in sentiment analysis

Figure 5 (a) demonstrates that there are 21 diaries devoted to exploring feeling examination about Twitter information. Utilization of semantics (Figure 1 (b)) is likewise observed to be higher. Nonetheless, there is less spotlight on philosophy based approach as appeared in Figure 1(c), while approaches that consolidate utilizes semantics and metaphysics is to a great degree less to discover in the current framework. General perception is that examination towards notion investigation is seen with less number of research papers when contrasted with inquire about issues in different surges of information examination or mining approaches. Figure 1 demonstrates the much of the time tended to issues in the zone of notion examination. The following area talks about the ongoing strategies of feeling investigation. Figure 2,3,4 shows the steps of data extraction,processing,xlsx file generationtraning phase.

VI. CONCLUSION

The investigation result straightforwardly demonstrates that current research systems have less underscored on the new types of the dynamic information that is described by higher multifaceted nature level as structure, heterogeneity, vulnerability, and so forth., There are different noteworthy research issues e.g. computational proficiency, unstructured information, selection of complex instances of open audits, and so forth, which is gotten a less consideration from the exploration network. Concentrates towards for bigger and unstructured information will require more consideration in the territory of supposition examination to guarantee the selection of conclusion investigation for up and coming correspondence related advancements. Accordingly, our future work will be toward tending to such open research issues. We will start our examination utilizing unstructured information and acquaint a system with countermeasure the multifaceted nature engaged with supposition mining from such complex information

REFERENCES

- [1] . HaCohen-Kerner, Y., & Badash, H. (2016). Positive and Negative Sentiment Words in a Blog Corpus Written in Hebrew. *Procedia Computer Science*, 96, 733-743.
- [2] Vidya, N. A., Fanany, M. I., & Budi, I. (2015). Twitter Sentiment to Analyze Net Brand Reputation of M obile Phone roviders. *Procedia Computer Science*, 72, 519-526
- [3] Kolchyna, O., Souza, T. T., Treleaven, P., & Aste, T. (2015).Twitter Sentiment Analysis: Lexicon Method, Machine Learning Method and Their Combination. arXiv preprint arXiv:1507.00955.
- [4] Amarouche, K., Benbrahim, H., & Kassou, I. (2015). Product Opinion Mining for Competitive Intelligence. *Procedia Computer Science*, 73, 358-365.
- [5] Kolchyna, O., Souza,T. T., Treleaven, P., & Aste, T. (2015).Twitter Sentiment Analysis: Lexicon Method, Machine Learning Method and Their Combination. arXiv preprint arXiv:1507.00955.

- [6] Medhat, W., Hassan, A., & Korashy, H. (2014). Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal*, 5(4), 1093-1113
- [7] Musto, C., Semeraro, G., & Polignano, M. (2014). A comparison of Lexicon-based approaches for Sentiment Analysis of microblog posts. *Information Filtering and Retrieval*, 59
- [8] Basari, A. S. H., Hussin, B., Ananta, I. G. P., & Zeniarja, J.(2013).Opinion mining of movie review using hybrid method of support vector machine and particle swarm . *Procedia optimization Engineering*, 53, 453-462..
- [9] Pak, A., & Paroubek, P. (2010, May). Twitter as a Corpus for Sentiment Analysis and Opinion Mining. In *LREc (Vol.10, pp.1320-1326)*.
- [10] Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and trends in information retrieval*, 2(1-2),1-135.
- [11] Hu, M., & Liu, B. (2004, August). Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 168-177)*.ACM).
- [12] Thompson, A. R., & O'Loughlin, V. D. (2015). The Blooming Anatomy Tool (BAT): A discipline-specific rubric for utilizing Bloom's taxonomy in the design and evaluation of assessments in the anatomical sciences. *Anatomical sciences education*, 8(6), 493-501
- [13] Bhargav H S, Application of Blooms Taxonomy in day-to-dayExaminationsIEEE(2016)
- [14] Chowdhry, B. S. (2013). Successful transformation of ICT graduate program: A role model for developing countries. *Wireless personal communications*, 69(3), 1013-1023.