# MACHINE LEARNING IN A TWITTER SENTIMENT ANALYSIS OF HEALTH CARE TWEET & USING BIG DATA FRAME WORK

J.Uma, Research scholar
Department of computer science
Periyar University PG Extensions centre, Dharmapuri.

Dr.K.Prabha, Assistant Professor
Department of computer science
Periyar University PG Extension centre, Dharmapuri.

## Abstract

Content mining has moved toward becoming pole mainstream inquire about zone. It can manage AI utilizing content investigation. It contains unstructured content which enormous measure of data can't just strategy by preparing PC and information from unstructed content finished by content mining. Content mining manages numerous strategies just as from data recovery, data extraction and furthermore common language handling and associate with calculation and techniques for KDD, web index and data reterival have most explicit pursuit Query language. This exploration fields use information mining calculation, for example, grouping, bunching, associatioie rule. This paper contains, though commentator of content mining strategy devices, and different applications.

Key words: Text mining, Data mining, Natural language process, Machine learning, social media.

## INTRODUCTION

Content mining is characterized as "The division of concealed and possibly needful data from literary information [1]. Content mining as new pursuit to extricate significant data from content language as normal. The present culture, content is most business correspondence formal trade of data. Content mining is same as information mining with the exception of information mining instruments [2] are planned unstructured, or semi organized information like as messages, content information base, picture, sound and video and so on. Content mining is a superior arrangement. Content mining is the procedure of obtained learning from printed database [3]. Content mining is a multi disciplinary field dependent on data extraction, data reterival, information mining, Machine Learning, Natural language preparing, and aggregation semantics [4], It can play with semi-organized and unstructured information, formal. Content mining systems are connected with industry, scholastic, web improvement zone, and different fields [5]. Content mining applications, for example, web crawler, Customer relationship the board, extricate data, Predictive investigation, misrepresentation identification and online networking examination, use content digging

for supposition mining, highlight extractor, estimation, prescient, and pattern examination [6]. Various practices of sharing and conveying are not based the substance but rather likewise based on reiteration of the content4. In the ongoing period miniaturized scale blogging has turned out to be very common21 and mainstream stage for every single online client. Millions/Billions of clients are imparting their insight on different angles on well known and in vogue sites, for example, twitter, Face book, tumbler, flash, LinkedIn etc.5 Twitter is a celebrated miniaturized scale blogging and person to person communication administration which gives the office to clients to share, convey and translate 140 words' post known as tweet3, 6. Twitter have 320M month to month dynamic client. Twitter is open through site interface, SMS, or cell phones. 80% clients are dynamic through mobiles7. In the miniaturized scale blogging administrations clients commit spelling errors, and use emojis for communicating their perspectives and emotions13. Characteristic language handling is additionally assuming a major job and can be utilized by the feelings communicated.

### I.LITERATURE SURVEY

A few calculations have been connected in the field of nostalgic examination in the course of recent years [13]. Dwindle D.

Turney [9] proposed a solo methodology of SO-PMIIR to order surveys as approval (positive) or disapproval (negative). The watched precision shifted from 66% for film surveys to 84% for vehicle audits. V K Singh et al [3] likewise utilized SO-PMI-IR to mine the understudies' conclusion in regards to

various subjects by gathering input from them in printed design.

P Walia et al [1] investigated solo (SO-PMI-IR) just as administered approach (NB and SVM) for nostalgic examination of film audits. The outcomes demonstrated that SOPMI-IR gave the best exactness and NB beat SVM. V K Singh et al [2] likewise investigated SWN strategy alongside NB, SVM, thus PMI-IR on film surveys.

Xing Fang and Justin Zhan [5] proposed another element vector age calculation to perform feeling extremity arrangement of item surveys (magnificence, books, home, and gadgets) acquired from amazon.com. T. K. Shivaprasad and J. Shetty [4] displayed the scientific categorization of different wistful investigation calculations. They investigated NB, SVM and ME put together managed methodologies with respect to surveys from games, gadgets, and PC.

Rodrigo Moraes et al [8] considered the presentation of SVM and fake neural system (ANN) on book dataset and referenced the extent of wistful examination in the book space.

## II.RELATED WORK

The related work talked about above has been condensed in table I. As clear from the table, NB, SVM, thus PMIIR are the most encouraging strategies in the field of wistful investigation. principle work in conclusion investigation just as sentiment mining, for example, impartial, negative and positive polarities in the assessment [14]. For instance, of positive assessment, "Coke tastes great" demonstrates that coke is delectable, a positive statement.However, "Coke beverages are not beneficial whenever expended in huge amount" shows that Coke beverages are destructive to body on the off chance that they drink a ton, a negative articulation. With respect to the impartial explanation, "Coke is a carbonated beverage". This paper clarifies the slant investigation scientific categorization or systems dependent on the data from [11][12][13][16]. In the event that you wish to have point by point proposed assumption investigation taxonomy, refer to this paper [15].

The conclusion can be classified as normal and near feeling [17]. Case of normal feeling, "The showcase quality is fresh." demonstrates that the part of "picture quality" is alluded legitimately and it gives a positive extremity [17]. With respect to the aberrant customary opinion," After applying the cream, my skin broke out completely." indicates that "the cream" in a roundabout way express the cream is awful for the skin and consequently it gives a negative extremity [17].

With respect to the similar conclusion, "The processor speed and screen goals of are superior to IPhone 6; anyway the metal group of iPhone 6 is more alluring than S6."indicates S6 has two positive feeling and one negative supposition [17]. It verifiably express the processor One Plus is superior to Yureka [17] .

Because of unpredictable, short type of content short length and slang content of tweets it is trying to foresee extremity of notion content. In estimation a blend of utilizations are expected to ponder and these all requests enormous number of conclusions from opinion holder. A rundown of supposition is required, as in extremity disambiguation and investigation; a solitary assessment isn't sufficient for choice. A typical type of notion investigation is viewpoint based for example telephone,

quality, voice, battery and so on. Rafael Michal Karampatsis8 et al. depicted the twitter slant examination for determining the extremity of messages. They utilized the two phase pipeline approach for investigation. Creators utilized the entirety classifier at each stage and a few highlights like morphological, POS labeling, vocabulary and so forth are distinguished. Joao Leal et al.11 attempted to arrange extremity of messages by utilizing AI draws near. Joachim Wagner et al. depicted work on perspective based extremity characterization by utilizing managed AI with Lucie Flekova et al.10 additionally took a shot at notion extremity expectation in twitter text.Nathon Aston et al.3 chipped away at estimation examination on OSN. They utilized a stream calculation utilizing changed adjusted for assumption examination. Lifna C.S.4 advances a novel methodology where the different themes are gathered into classes and afterward dole out weight age for each class by utilizing sliding window preparing model upon twitter streams. In the comparative way Emma Haddi et al.12 examined the job of content pre-handling for estimation analysis.EfthymiosKouloumpis14 characterized and clarified three way feeling investigation in twitter for recognize positive, negative and unbiased assessments. Efstratios Kontopoulos16 proposed a novel methodology for examination of opinion. The methodology is philosophy based and it just discover the estimation score just as evaluation for each particular thought in the post.

Semantic examination is the examination of individuals' suppositions, convictions, dispositions and feelings towards a substance, for example, items, administrations, occasions, issues and subjects [1]. It is the field of AI which has picked up the consideration of analysts since the start of the century. Mill operator et al. [12] presents WordNet, an online database for English language semantic handling utilizing equivalent word sets (synsets) relationship. SentiWordNet [13] is a progression of WordNet as an apparatus for learning based word level handling by means of structure a lexicon to discover a score of each word.

Kim and Hovy [16] worked on a word granularity by utilizing at first some seed words and utilizing them to make a net; they continued further to sentence level by consolidating the qualities of the words, as they arrange individuals' assessments. Additionally, Wilson et al. [17] worked on an expression level, by running a regulated learning way to deal with decide the extremity or nonpartisanship of expressions. Moreover, archive granularity [18] utilized word recurrence and grammatical feature approach on Amazon audits in classifications, similar to books, DVDs, electronic and kitchen apparatuses to assess the reaction of individuals about the items.

Twitter spilling API1 was utilized to assemble information for item assessment investigation [3]. The point of utilizing twitter information is to comprehend general supposition. Around 60,000 tweets were gathered utilizing Twitter API to investigate client feelings on broadly utilized cell phones in Korea [21]. Kumar et al. exhumed assessments of the individuals about the nature of administrations given via Airtel organization [22]. For this reason, they gathered 80,000 tweets utilizing the hashtag "#Airtel". They surveyed them utilizing Naïve Bayes approach with a precision of 80.9% on Mahout introduced over Hadoop to characterize them into various classes. They utilized term recurrence and backwards report recurrence for inside preparing

## III.SENTIMENTAL ANALYSIS IN CATEGORIES

**Lexical Analysis**: it expects to ascertain the extremity of an archive from the semantic direction of words or expressions inside the reports. Nevertheless, an application alludes to vocabulary examination and it doesn't reflect to think about the contemplated setting.

**Machine Learning**: it incorporates building models got from named preparing dataset (sentences or occasions of writings) so as to discover the archive direction. Concentrates that apply to this kind of methods have been executed on a definite subject.

The utilization of normal language preparing (NLP) by assumption investigation just as feeling mining is to gather and look at the estimation words and sentiments [7]. Along these lines, finding abstract demeanors in the huge social information is viewed as a well known territory in the field of NLP and information mining [8]. Feeling examination helps to accomplish various objectives like watching open disposition with respect to showcase

## IV. CLASSIFICATION TECHNIQUE

Characterization is a method which serves to completely characterize in which informational collection completes a specific information occasion fall into. In content mining, all the content classifiers can work on an enormous measure of information as indicated by their separate imperatives. In K-closest neighbor classifier classes may not be fundamentally required to be straightly discernable but rather in this classifier, it is truly tedious to discover the closest neighbors if information is tremendous. In SVM classifier, the exactness of results can be high, however it is intricate and requires more existence in both preparing and testing [3]. In ANN classifier, it works very well with just a couple of parameters to modify, yet the preparing time can be truly elevated if the neural system is huge. In this paper, probabilistic Naïve Bayes classifier has been utilized whose usage isn't just basic, yet in addition has astounding effectiveness and arrangement rate [7,8]. Likewise, according to the information size taken in this paper, this calculation demonstrates to give the best outcomes and consequently being most suitable in content arrangement of the information.

## V.METHODOLOGY

### 5.1. Twitter Data in Python

Twitter Data in Python. Twitter information is in contrast with the data shared by the greater part of the other interpersonal interaction destinations since it reflects information that the clients pick to share transparently out in the open. The twitter API stage gives extensive access to open tweets that clients over the world have bestowed. So as to get to the twitter API, the accompanying methodology has been adopted. The preeminent advance to bring the tweets from twitter has been to make a twitter application to gain admittance to the Twitter engineer account by the indistinguishable username as the one signed into. This has been done so as to acquire the certifications that are expected to stream the tweets from the twitter API. Further, utilizing a python library called Tweepy the tweets were brought from the twitter API [9]. Tweepy empowers python to cooperate with twitter API and thus gushing of the tweets from the twitter. The tweets so acquired have been coordinated into a

json document. In this paper, the system has been executed by getting tweets by utilizing catchphrase Kashmir and consequently an information ordinarily of size 339MB has been gotten..

### 5.2 Preprocessing

Content mining new field to extricate significant data from characteristic language. It is the procedure of explicit concentrate data to individuals. Content mining manages content is the most common route for the formal trade and supposition. The most significant of procedure of content mining advancement in measurable, scientific, Linguistic and Pattern acknowledgment strategies, it permit programmed investigation of unstructed data just as concentrate that data as opposed to looking through words,, The overlay looking at higher level.

**Process:**

Text mining involves a series of communication and mine information,

**Text Preprocessing:**

1. Text Data base
2. Text Preprocessing

   * To kenization

   * Stop word removal

   * Stemming

3. Text transformation

   * Feature Generation

4. Feature selection

   * Attribute selection

5. Text Mining Techniques

6. Evaluation

**Text Cleanup:**

# removing of any unnecessary or unwanted information removes ads from web Pages.

Tokenization:

# Split the text on white Space and at punctuations

Part- of – Speech:

Word class assignment to each to token it is input given by tokenized text

**Text transformation**:

It is also known as variable selection the main assumption when using future selection technique, The data contain may redundant or irrelevant futures, Future selection is a dividing more general of future selection.

**Data mining:**

Text mining is process join with the data mining process, Evaluate to check the result of correctness.

**v. Sentimental Analysis using Different Algorithm**

In the AI field, arrangement techniques have been created, which utilize various methodologies to order unlabeled information. Classifiers could require preparing information. Instances of AI classifiers are Naive Bayes, Maximum Entropy and Support Vector Machine [14] [15, 16]. These are ordered as managed AI strategies as these require preparing information. Mention that preparation a classifier successfully will make future forecasts simpler. AI approach is utilized to prepare a calculation with a predefined dataset before applying it to genuine dataset. AI methods first trains the calculation with some specific contributions with known yields so later it can work with new obscure information. The absolute most prestigious works dependent on AI are as per the following.

**Support  Vector Machines (SVM)**

A standard SVM takes an accumulation of enormous information and predicts, for each given contribution, there are some achievable classes which structures the yield. At the point when given an accumulation of preparing models, each set apart as having a place with a chosen class, a SVM preparing guideline manufactures a model which will be utilized to relegate new models into a class [8]. A SVM model might be a portrayal of the models as focuses in territory, mapped, for example, the individuals from the different classes are separated by a hole as wide as feasible. New models are then mapped into that extremely same territory and expected to have a place with in any event one of the classes bolstered that part of the hole they fall in. Characterizing in all respects officially, a help vector machine builds a hyper plane or an accumulation of hyper planes in an unending dimensional zone, which might be utilized for grouping. Normally, a compelling partition is accomplished by the hyper plane that has the most significant separation to the nearest preparing data of any classification. Bigger the edge, lower would the speculation mistake of the classifier be[9].

**Naïve Bayes**

This methodology assumes the supply of at any rate a lot of articles with pre-relegated sentiment and reality names at the record level [10]. They utilized single words, while not stemming or stop word expulsion as choices. Credulous Bayes relegates an archive d to the classification c, that amplifies P (c/d) by applying Bayes' standard.

## VI.SENTIMENT ANALYSIS AND OPINION MINING FRAME WORK

We present a framework for sentiment analysis which includes data collection, pre-processing, sentiment score calculation for tweets, classification and polarity prediction.

$$\mathbf{P}(C|X) = P \, x \, c \, P(c)/P(x) \text{-------------------(1)}$$

Here P (c | x) is a Posterior probability.
P (x | c) is a Like hood.
P(c) is a Prior post probability.

P(X) is a predicator probability

Step1: It starts with frequency table.
Step2: Create Like hood table by finding the probabilities overcast probabilities.
Step3: Applying that probability theorem.

Step4: Stop the process.

Table1: Collect from the table

| Patients | Affect |
|----------|--------|
| Flu | Yes |
| Dengue | Yes |
| Ebola | No |
| Dengue | Yes |
| Flu | Yes |
| Ebola | No |
| Flu | yes |
| Ebola | No |
| Flu | yes |
| Ebola | No |

**Table2:Frequency table**

| Affect | NO | Yes |
|--------|-----|-----|
| Flu | - | 4 |
| Dengue | 3 | 2 |
| Ebola | 2 | 3 |
| Total | 5 | 9 |

**Table3:Result**

| Affect | No | Yes | |
|--------|-----|-----|---|
| Flu | - | 4 | 4/14=0.29 |
| Dengue | 3 | 2 | 5/14=0.36 |
| Ebola | 2 | 3 | 5/14=0.36 |
| Total | 5/14=0.36 | 9/14=0.64 | |

(**yes** Ebola) = p(Ebola  **yes**) * p (**yes**) / p (Ebola )

P (Ebola | yes) = 3/9 = 0.33 p (Ebola )
=5/14=0.36,p (yes) = 9/14 = 0.64

**Now p (yes | Ebola) = 0.33 * 64 / 0.36 = 0.60**

## VII. Conclusion

There are unmistakable Emblematic and AI systems to perceive thoughts from substance. Simulated intelligence techniques are less perplexing and capable than Representative strategies. These methodologies can be associated for twitter thought examination. There are certain issues while overseeing perceiving eager catchphrase from tweets having various watchwords. It is moreover difficult to manage erroneous spellings and slang words. To deal with these issues, a beneficial segment vector is made by doing feature extraction in two phases after genuine preprocessing. In the underlying advance, twitter unequivocal features are removed and included to the component vector. Starting now and into the foreseeable future, these features are emptied from tweets and again feature extraction is done as if it is done on run of the mill content. These features are in like manner added to the component vector. Course of action precision of the component vector is taken a stab at using particular classifiers like Nave Bayes, SVM, Greatest Entropy and Gathering classifiers. All of these classifiers has for all intents and purposes equivalent exactness for the new segment vector. This component vector performs well for electronic things.

## REFERENCES

1Neha S. Joshi, Suhasini A. Itkat, " A Survey on Feature Level Sentiment Analysis" (IJCSIT) International Journal of ComputerScience and Information Technologies, Vol. 5 (4) , 2014, 5422-5425.

2. He Y., "Incorporating sentiment prior knowledge for weakly supervised sentiment analysis ", ACM Transactions on Asian Language Information Processing, Vol. 11(2).

3. N. Veeranjaneyulu, Akkineni Raghunath, B, Jyostna Devi, Venkata Naresh Mandhala, "Scene Classification Using Support Vector Machines With Lda " journal of theoretical and applied information technology 31 may 2014. Vol. 63 No.3

4. Ankush Sharma, Aakanksha, "A Comparative Study of Sentiments Analysis Using Rule Based and Support Vector Machine", IJRCCE Vol. 3, Issue 3, March 2014.

5. P. Saloun, M. Hruzk and I. Zelinka, "Sentiment Analysis e-Bussines an e- Learning Common Issue," ICETA 2013 ,11th IEEE International Conference on Emerging eLearning Technologies and Applications, Stary Smokovec, The High Tatras, Slovakia, October
24-25, 2013.

6. A. Tamilselvi, M. ParveenTaj, "Sentiment Analysis of Micro blogs using Opinion Mining Classification Algorithm " International Journal of Science and Research (IJSR) ISSN (Online): 2319-7064 Volume 2 Issue 10, October 2013.

7. Pablo Gamallo, Marcos Garcia, "Citius: A Naive-Bayes Strategy for Sentiment Analysis on English Tweets" Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014), pages 171–175, Dublin, Ireland, August 23-24 2014.

8. Ankitha Srivastava, Dr.M.P. Singh, "Supervised SA of product reviews usin Weighted k-NN Algorithm," 2014 11th International Conference on Information Technology.

9. Kerstin Denecke, "Using SentiWordNet for Multilingual Sentiment Analysis," ICDE Workshop 2008, 978-1-4244-2162-6/08/ 2008 IEEE.

10. Brett W. Bader, W. Philip Kegelmeyer, and Peter A. Chew "Multilingual Sentiment Analysis Using Latent Semantic Indexing and Machine Learning," 2011 11th IEEE International Conference on Data Mining Workshops.

11. Lizhen Liu, Xinhui Nie, Hanshi Wang, "Toward a Fuzzy Domain Sentiment Ontology Tree for Sentiment Analysis," 5th Internationa Congress on Image and Signal Processing (CISP 2012).) 2012.

12. Swati A. Kawathekar1, Dr. Manali M. Kshirsagar, "Movie Review analysis using Rule-Based &Support Vector Machines methods",IOSR Journal of Engineering Mar. 2012, Vol. 2(3), March. 2012, pp: 389-391.

13. Blinov P. D. ,Klekovkina M. V. , Kotelnikov E. V. , Pestov O. A." Research of lexical approach and machine learning methods for sentiment analysis", Vyatka State Humanities University, Kirov, Russia.

14. Klekovkina M. V., Kotelnikov E. V., "The automatic sentiment text classification method based on emotional vocabulary" , Digital libraries: advanced methods and technologies, digital collections (RCDL-2012) , pp. 118–123.

15. Lan M., Tan C. L., Su J., Lu Y. (2009), "Supervised and traditional term weighting methods for automatic text categorization", IEEETransactions on Pattern Analysis and Machine Intelligence, Vol. 31(4), pp. 721–735.

16.Hatzivassiloglou, V., Wiebe, J.*Effects of Adjective Orientation and Gradability on Sentence Subjectivity*. Proceedings of the 18th International Conference on Computational Linguistics, New Brunswick, NJ. 2000.

17. Kennedy, A., Inkpen, D.*Sentiment Classification of Movie and Product Reviews Using Contextual Valence Shifters*. Computational Intelligence.2006. pp. 110-125.

18. Turney, P. *Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews.* Proceedings of ACL,
Philadelphia, PA. July 2002. pp. 417-424.

19. Kamps, J., Marx, M., Mokken, R. J.*Using WordNet to Measure Semantic Orientation of Adjectives*. LREC 2004. Volume IV, pp. 1115-1118.

20. Andreevskaia, A., Bergler, S., Urseanu, M.*All Blogs Are Not Made Equal: Exploring Genre Di_erences in Sentiment Tagging of Blogs.*International Conference on Weblogs and Social Media (ICWSM-2007), Boulder, CO. 2007.

21. Hemnaath, R., and Low, B. W. 2010. *Sentiment Analysis Using Maximum Entropy and Support Vector Machine*. Semantic Technology and Knowledge Engineering in 2010. Kuching,Sarawak

22. Turney, P.D., Littman, M.L.*Measuring Praise and Criticism: Inference of Semantic Orientation from Association.* ACM Transactions on Information Systems. 2003. pp. 315-346.

23. B. Pang, L. Lee and S. Vaithyanathan. *Thumbs up?: sentiment classification using machine learning techniques.* In EMNLP '02: Proc. of the
ACL-02 conf. on Empirical methods in natural language processing, pages 79–86. ACL, 2002.

24. Das and Chen.*Yahoo! for Amazon: Sentiment Extraction from Small Talk on the Web Management Science* 53(9), pp. 1375–1388, © 2007