

e-ISSN: 2395 - 7639



INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH

IN SCIENCE, ENGINEERING, TECHNOLOGY AND MANAGEMENT

Volume 10, Issue 2, February 2023



INTERNATIONAL STANDARD SERIAL NUMBER INDIA

 \mathbb{X}

Impact Factor: 7.580

ili in the second secon

| ISSN: 2395-7639 | www.ijmrsetm.com | Impact Factor: 7.580 | A Monthly Double-Blind Peer Reviewed Journal |

Volume 10, Issue 2, February 2023

| DOI: 10.15680/IJMRSETM.2023.1002009 |

Multi-Task Ensemble Model for Aspect and Polarity Classification

C.A.KANDASAMY, JAYABALAAGI K S

Assistant Professor, Department of Computer Applications (MCA), K.S.R. College of Engineering (Autonomous),

Tiruchengode, India

Department of Computer Applications (MCA), K.S.R. College of Engineering (Autonomous), Tiruchengode, India

ABSTRACT: Item detection from a tweet is a common task to understand the current movies/topics attracting a large number of common users. However the unique characteristics of tweets (short and noisy content, and a large data volume) make the item detection a challenging task. Existing techniques proposed for item detection uses battery of one class classifier using key word matching techniques and SVM classifier and those techniques provide better accuracy but the features are extracted are found to be noisy, this is a major limitation in SVM classifier.

In this system a SVM classifier with genetic algorithm optimization is proposed. In GA optimization we use 'accuracy 'of SVM as a fitness function; only the best features are selected. And this will improve accuracy for item detection and also the system provides user rating based on the polarity of tweets. This system is expected to improve in terms of classification accuracy when GA is combined with SVM

KEYWORDS: Item detection, Polarity detection, Twitter, SVM, GA

I.INTRODUCTION

Posting comments about TV programs/Movies using second screen devices (e.g. tablets) is very common. The **Twitter** is the most known and used micro blogs among the preferred social services to post short messages while watching TV. In order to elicit user preferences from a micro blog comment, we need to recognize if the user writing about specific TV shows or Movies .However automatically detecting the subjects of the tweets is a challenging task because of 140 –character limit of tweets (comments) strongly affects the way people write on twitter. Tweets are unstructured data. Tweet involve

- Minimal contextualization
- Use of slangs
- Abbreviations
- Tiny URLs, etc., so the tweet has more noise and difficult to understand.

An example tweet:

Anyone going to @cstn this year? @theperfectfoil will be presenting BLJ in Imaginarium tent onThurs evening, July 5th

Refers to a talk that would be given by Steve Taylor ('The Perfect Foil') about the movie blue like jazz ('BLJ') during Cornerstone music festival in Illinois ('cstn') on Thursday ('Thurs').

In this work a solution to analyze the content of tweets and to identify whether they refer to one (or more) known items is proposed. e.g., movies or TV programs (i.e., Item detection) and item detection prevents the use of traditional classifiers because in that if they want to add to/remove the item means the classifiers required to retrained instead of that here battery of one-class classifiers are used ,each one class classifier is trained for each item in the catalog. When a new item (i.e., class) is added to the catalog, we simply need to train an additional one-class classifier, while no updates are required for the remaining one-class classifiers previously trained.

Implementation of each one-class classifier as a pipeline composed by three stages is decided, An unknown tweet is initially processed by stage 1, that can either assign the tweet to its associated class or it can discard the tweet. Only tweets not classified at stage 1 are processed by stage 2. Similarly, only tweets not classified by stage 2 are processed by the last stage. Tweets not even classified by stage 3 remain not classified by this specific classifier. We used two keyword matching algorithms in stages 1 and 2: these algorithms are based on regular expressions. As for stage 3, we



| ISSN: 2395-7639 | www.ijmrsetm.com | Impact Factor: 7.580 | A Monthly Double-Blind Peer Reviewed Journal |

Volume 10, Issue 2, February 2023

DOI: 10.15680/IJMRSETM.2023.1002009

implemented a more advanced approach based on the support vector machine (SVM) and genetic algorithm (GA) algorithm for item detection. After that we can identify the polarity of the tweets of that item using GA based on the positive and negative tweets aggregate counting providing star rating for user reference.

II. EXISTING SYSTEM

An unknown tweet is initially processed by stage 1, that can either assign the tweet to its associated class or it can discard the tweet. Only tweets not classified at stage 1 are processed by stage 2. Similarly, only tweets not classified by stage 2 are processed by the last stage. Tweets not even classified by stage 3 remain not classified by this specific classifier. We used two keyword matching algorithms in stages 1 and 2 are based on pattern matching. As for stage 3, we implemented a more advanced approach based on the SVM algorithm



2.1 Stage 1: Exact Match

This stage deploys a simple keyword matching on the basis of the movie title, denoted as exact match. Given the item-i classifier, this algorithm classifies the tweet as belonging to class i if the text contains exactly the title of movie i.

E.g. Tweet: I like Harry Potter

Item name: Harry Potter

2.2 Stage 2: Free Match

Stage 2 of the one-class classifiers is composed of a keyword matching algorithm, denoted as free match. Differently from the keyword matching implemented in stage 1, stage2

(i) Searches for single words of the title within the tweet and

(ii) Searches for words of the title contained in other words

In the tweet. For instance, this stage is able to find a matching between the tweets

E.g. the future is back to the present in the past

And the movie **`Back to the Future'**. Still, this approach allows us to match the movie title also with hash tags and mentions. For instance, the tweet

E.g. atabduction just got home and it was awesome!!!!

Is matched by stage 2 because it contains the mention `atabduction', but not by stage 1.

2.3 Stage 3: SVM-Based

2.3.1 Text Preprocessing Tasks

Tokenization, Stop Word Removal, Stemming is performed.

2.3.2 Feature Selection

The following four types of features are extracted:



| ISSN: 2395-7639 | www.ijmrsetm.com | Impact Factor: 7.580 | A Monthly Double-Blind Peer Reviewed Journal |

Volume 10, Issue 2, February 2023

| DOI: 10.15680/IJMRSETM.2023.1002009 |

- **Unigrams** correspond to the single terms. Each distinct unigram is represented by a dimension in the feature space. Unigrams formed by stop words are discarded.
- **Bigrams** are pairs of adjacent terms. Each unique bigram corresponds to a dimension in the feature space.
- **Hash tags** represent the tags given by a user to a tweet. However, due to the lack of constraints and common criteria, users can freely use hash tags, either reusing existing tags or defining new ones. Any different hash tag is represented with a dimension in the feature space.
- **Titles** are not linked to specific terms, but it refers to the degree of matching between the tweet content and the item title. Each possible title related to an item is represented by a dimension in the feature space.

2.3.3 Feature Normalization

Binary weighting, The TF-IDF schema are used for normalizing the tweets

2.3.4 Classification

Differently from the previous two stages, the SVM algorithm requires also a learning phase to tune its parameter for a specific class. Learning is performed using a 5-fold Cross validation on a set of classified tweets. The output of the SVM algorithm is a real number, referred to as decision value. Only tweets with a decision value over a fixed threshold are classified.

The performances of the classification system have been measured using a standard hold-out dataset partitioning. For each one-class classifier, we randomly selected 200 tweets to form the test set. The class (i.e., the item) of tweets in the test set is assumed to be unknown to the classifier and is used to verify if the predicted class corresponds to the actual class. Part of the remaining tweets has been used to form the training set, used by stage 3 as sample tweets for the learning phase.

2.4 Drawbacks of Existing System

- SVM-Features are selected with noisy data.
 - -Features with noisy information increase the complexity of classification.
- Rating of the item is not specified

III. PROPOSED SYSTEM

Stage 1 and stage 2 are similar to existing system but in stage 3 for movie detection from tweets feature selection will be done by GA and classification will be done by SVM to improve the item detection accuracy. After detecting the item take all the tweets of that item detect the polarity using sentiment analysis by EWGA. After detecting polarity provide user rating for user reference based on aggregate counting of positive and negative tweets of that item.

3.1 Item Detection by GA+SVM

In stage 3 after text preprocessing task feature extraction will be done by GA is represented in the figure 4.2, the same four existing features like unigrams ,bigrams ,hash tags and titles are extracted at last the best feature will be obtained.



ili in the second secon

| ISSN: 2395-7639 | www.ijmrsetm.com | Impact Factor: 7.580 | A Monthly Double-Blind Peer Reviewed Journal |

Volume 10, Issue 2, February 2023

| DOI: 10.15680/IJMRSETM.2023.1002009 |

3.2 Polarity Detection By EWGA

After detecting the item name we will collect all tweets of that item then we will do the EWGA classification for polarity detection. For this purpose we will extract following type of features

- Parts of speech
- Senti word net
- Frequency
- Stemming
- Chunk label
- Dependency parser
- Positional aspect
- Term distribution
- Thematic word

After applying GA the best feature is selected as a subjective word based on that identifies its polarity whether it is positive or negative.

3.3 User Rating Analysis

After detecting the polarity of user reviews count the number of positive tweets and negative tweets based on the count provide user rating for the movies in the form of stars specified in the below Figure.



3.4 Advantages of Proposed System

- GA+SVM will improve the accuracy.
- Eliminate noise during feature extraction
- EWGA will identify polarity of the tweets.

IV. CONCLUSION

In this work

- Best feature will obtained using GA, so the performance of SVM classifier will increase for item detection.
- User rating of the movie will be obtained from polarity detection using GA

REFERENCES

- [1] T. Joachim's.. "Text categorization with support vector machines: Learning with many relevant features" In Proceedings of the 10th European Conference on Machine Learning, ECML '98.
- [2] M. Manevitz, M.Yousef, "One-Class SVMs for Document Classification", Journal of Machine Learning Research 2 (2001) 139-154 Submitted3/01; Published12/01.
- [3] X.Ming Zhao, De-S. Huang, Y.Cheung, H.Wang and X.Huang "A novel hybrid GA/SVM system for protein sequences classification" IDEAL 2004, LNCS 3177, pp. 11–16, 2004. Springer-Verlag Berlin Heidelberg 2004
- [4] Liaoyang LIU, Hui FU "A Hybrid Algorithm for Text Classification Problem", PRZEGLĄD ELEKTROTECHNICZNY (Electrical Review), ISSN 0033-2097, R. 88 NR 1b/2012.
- [5] A.Das, S.B.Opadhyay"Subjectivity Detection using Genetic Algorithm", the 1st workshop on computational Approach, 2010.

IJMRSETM©2023



| ISSN: 2395-7639 | www.ijmrsetm.com | Impact Factor: 7.580 | A Monthly Double-Blind Peer Reviewed Journal |

Volume 10, Issue 2, February 2023

| DOI: 10.15680/IJMRSETM.2023.1002009 |

- [6] A.Abbasi,H.Chen,andA.Salem "Sentiment Analysis in Multiple Languages: Feature Selection for Opinion Classification in Web Forums", ACM Transactions on Information Systems, Vol. 26, No. 3, Article 12, Publication date: Jun 2008.
- [7] P.Cremonesi, R.Pagano, S.Pasquali, and R.Turrin, "TV Program Detection in Tweets," Proc. ACM, EuroITV'13, June 24–26, 2013.









INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY RESEARCH

IN SCIENCE, ENGINEERING, TECHNOLOGY AND MANAGEMENT



+91 99405 72462

🕥 +91 63819 07438 🔀 ijmrsetm@gmail.com

www.ijmrsetm.com