



A Novel Text Representation Method for Irony and Stereotype Spreaders Detection

¹K. Dinesh Murthy,²S. V. S. Chandravadan,³Karunakar Kavuri, ⁴Archana Gelli

¹B. Tech Student,²B. Tech Student,³Associate Professor, ⁴Assistant Professor

^{1,3,4}Computer Science and Engineering, ²Electrical and Electronics Engineering

¹Swarnandhra Institute of Engineering and Technology, Narasapur, Andhra Pradesh, India

²Birla Institute of Technology & Sciences, Pilani, Hyderabad, Telangana, India

^{3,4}Swarnandhra College of Engineering and Technology, Narasapur, Andhra Pradesh, India.

Abstract : ThisThe digital information is spreading massively through social media platforms like Facebook, Twitter, Blogs, Reviews, and Instagram etc., which leads to extremely biased, rumors or false information being consumed and shared by Internet users every day. Disinformation or fake news including irony and stereotype information becomes a major problem in our present society. The irony and stereotype news influences the public health, economy and even politics. Knowingly or unknowingly, the people are spreading irony and stereotype information in the social media like Twitter and Facebook. There is a need of techniques to identify the irony and stereotype news spreaders in the social media through that an alert can send to the people in the community to check whether the message is received from real users or irony and stereotype spreaders. This task is analyzed by the PAN competition organizers and conducted a competition on irony and stereotype spreaders detection on Twitter dataset in 2022. Most of the researchers presented solutions for irony and stereotype spreaders detection based on machine learning techniques or deep learning techniques in the competition. In this work, we proposed a novel text representation method for irony and stereotype spreaders detection. In the proposed method, we represent the terms as m-dimensional vectors by computing the weight of a term in all documents of dataset. Each vector value in term representation is the weight of a term in a specific document. The documents are represented as vectors by aggregating the vectors of terms that are contained in that document. Each document of the dataset is represented as m-dimensional vectors. These document vectors are trained with two machine learning algorithms such as support vector machine and random forest for predicting the accuracy of irony and stereotype spreaders detection. The random forest classifier attained best accuracy of 97.86% for irony and stereotype spreaders detection when compared with state-of-the-art methods.

IndexTerms - Text Representation, Term Weight Measures, Machine Learning Algorithms, Irony and Stereotype Spreaders Detection.

I.INTRODUCTION

In these days, the people are heavily relied on social media for different types of news like latest updates on movies, famous people, politicians, trending stories etc. This was become advantage to several people to spread false information in the form of irony and stereotype messages in social media. Irony is a rhetorical figure that consists of saying the opposite of what is meant, using a tone, gesture or words that insinuate the interpretation that should be made. Stereotypes are often used, especially in discussions of controversial issues such as immigration, sexism, and misogyny. The flooding of irony and stereotype messages affects the stock markets, news system and opinions of people. The irony and stereotype news spreaders create more damage to the society when compared with irony and stereotype news creators. Most of the people are not having any idea that the news in social media platforms coming from genuine sources or not. Therefore, there is a need to have tools that are capable of determining when a user is employing sarcasm or irony to affect other people or groups of them, since the large amount of information that is generated daily makes manual control of it impossible.

The irony and stereotype spreaders detection is a type of author profiling task. The prediction of traits such as gender, age, occupation and origin of a person are topics that have been widely studied in the field of Author Profiling (AP) [1]. Recent research has been dedicated to detecting social behaviours and psychological characteristics of users by applying the techniques used in AP task [2]. During the last 3 years, PAN has launched several AP tasks dedicated to identifying users on Twitter who spread harmful content, as well as being able to identify those profiles that constitute chat bots, due to their participation in this activity. Precisely the task of the year 2019 was dedicated to differentiating bots from human profiles [3] and in 2020 it was aimed at identifying those users who shared fake news on the network [4]. At PAN 2021, last year, the AP task proposed was Profiling Hate Speech Spreaders on Twitter 2021 [5], whose main objective was to determine, from a set of 200 tweets per author, whether or not a user profile was a hate speech spreader. In 2022 [6], the AP task was focused on determining ironic profiles on Twitter, paying special emphasis on those authors who use irony to spread stereotypes. The goal of the task is to classify authors as ironic or not based on the number of tweets with ironic content. Therefore, given the Twitter authors along with their posts, the main objective will be to profile those authors that can be considered ironic.

In this work, we proposed a novel text representation method for irony and stereotype spreaders detection. Based on the analysis of the given dataset in the competition, we observed that the Irony and Not Irony classes of documents are differentiated by using

the content the authors used in their tweets. In the proposed method, we extracted all informative terms after cleaning the dataset by using pre-processing techniques. All extracted terms are represented as vectors. The number of dimensions in a vector is equivalent to number of documents in the dataset. In this work, the PAN competition training dataset contains 420 author files. So, each term is represented as a 420 dimensional vectors. Each vector value is a weight of a term in that specific document. These term vectors are used for representing the document vectors. Each document in the dataset is represented as a vector by aggregating all the term vectors those are contained in the document. The proposed document representation method utilizes all the information specified in that document. Two machine learning algorithms are used to generate the model by training these document vectors. The random forest classifier shows best performance for irony and stereotype spreaders detection.

II. LITERATURE SURVEY

Social networks have been playing an important role in the life of human beings for the last years, they have become a way to express and share information widely. In them, many people create harmful and offensive content towards others, such as irony, sarcasm and the use of stereotypes to refer to certain groups in society [7]. Because the information shared on the internet grows very fast, it is necessary to have systems that can automatically detect this unwanted behaviour on networks. Maria Fernanda Artigas-Herold et al., described [8] an approach to the Profiling Irony and Stereotype Spreaders on Twitter (IROSTEREO) task promoted by PAN CLEF 2022, where they want to identify profiles of users who post ironic content on Twitter. Their proposal is to build models based on n-grams of characters and words, as well as non-English words in combination with SVM and RF as classification algorithms, and obtains a majority vote of those with the best results for each representation. The proposed solution reached an accuracy of 91.67%.

Sabur Butt et al., described [9] the model submitted by the team CIC for "Profiling Irony and Stereotype Spreaders on Twitter (IROSTEREO)" at PAN 2022. Irony profiling can help in identifying stereotype spreaders and can enhance the understanding of author behaviours. They proposed a methodology focusing on feature engineering to classify irony for long texts based on multiple linguistic and emotion-based features. They also extensively discussed the shortcomings of the data and the proposed task to provide the future research direction. The paper reveals the impact of robust feature engineering with a machine learning approach on the long social media texts in the author profiles. Authors initiated a dimensionality reduction method and reduced the number of features to 100,000 most frequent features. The selected features were concatenated with features such as hate speech (aggressive, hateful, targeted), emotions (surprise, joy, sadness, fear, disgust, anger, other) and sentiments (positive, neutral, negative) that give another dimension to the model understanding. They used the XGB classifier for feature importance on a subset of the training set and later used SHAP to create shapely values for the selected linguistic features. The proposed method achieved an accuracy of 87.22% on the test set.

Daniele Croce et al., addressed [10] the problem of classifying whether a Twitter user has spreading Irony and Stereotype or not. Authors used a text vectorization layer to generate Bag-Of-Words sequences. Then, such sequences are passed to three different text classifiers such as Decision Tree, Convolutional Neural Network, and Naive Bayes. Based on the performance of these classifiers, they decided to select the final classifier as SVM. To test and validate their approach, they used the dataset provided for the author profiling task organized by PAN@CLEF 2022. Over several cross fold validations, the proposed approach was able to reach a maximum binary accuracy on the best validation split equal to 0.9474. On the test set provided for the shared task, proposed model is able to reach an accuracy of 0.9389.

Irony is a curious mode of communication in which the speaker says something that wants the audience to be interpreted oppositely. Its automatic detection is a very challenging task due to its complex interpretation, and it has a significant potential for various applications in text mining. Social Media platforms like Twitter offer a vital chance to analyse this literary technique since users frequently utilize it to give their opinions. José Antonio García-Díaz et al., designed [11] a contribution for the 2022 PAN's shared author profiling task and its subtask concerning Stereotype Stance Detection. The former consists in determining whether the authors spread irony and stereotypes and the latter is focused on identifying stereotypes that can hurt vulnerable groups. The organizers provide a set compiled from Twitter to carry out the task. In particular, they have proposed a supervised learning pipeline consisting of a combination of Deep Learning techniques that utilizes context and non-context embeddings to address the binary classification. The resulting system reaches promising results, achieving the fifth-best score in the main task with an accuracy of 96.67%.

Catherine Ikae proposed [12] a solution for solving the problem of profiling irony and stereotype spreaders on twitter using a random forest model with features obtained by using chi square feature scoring. The task is to determine whether the author of a given Twitter feed in English spreads irony and stereotypes. The training sample contains timelines of authors sharing irony and stereotypes towards, for instance, women or the LGTB community. Transforming this question into binary classification problem which requires us to classify authors as ironic or not. Evaluation with 300 chi2 features shows an overall performance of accuracy is 0.9722. They discovered that such a figure is still too huge after evaluating the terms found in feature sets with thousands of terms. They also found that only the top m terms (e.g., m = 300) depicting the highest discriminating powers selected with the chi2 and PMI feature selection techniques.

The PAN 22 Author Profiling Shared Task (IROSTEREO) aims to profile authors spreading irony and stereotypes on Twitter. Hyewon Jang experimented [13] with different classification methods such as traditional n-gram approach, state-of-the-art language models, and lexical approach using LIWC. As a baseline, they experimented with Term Frequency-Inverse Document Frequency features. After converting the data into a TF-IDF matrix using different N-gram sizes (N=1, 2), unigrams (N=1) proved to be the best in terms of classification performance and computational efficiency. In order to obtain results with better explainability, they also developed a model using only lexicon-based features. They used the software Lexical Inquiry and Word Count (LIWC) to extract features belonging to various lexical categories such as Psychological Processes and Linguistic Dimensions. A total of 93 features were extracted. They experimented with several traditional classifiers such as Random Forest Classifier (RF), Support Vector Classifier (SVC), Gaussian Process Classifier (GPC), Decision Tree Classifier (DT), and Adaptive Boost Classifier (ABC). They observed that the best result was obtained from the lexicon-based approach (LIWC) with the accuracy score of 0.88 on the validation data and 0.92 on the official test data.

Tiago Filipe Nunes Ribeiro et al., presented [14] a model for classifying irony and stereotype spreaders on Twitter based on the dataset provided for the PAN2022 task of IROSTEREO. They take a feature engineering approach focusing on lexical and

stylistic features and improve on the character n-gram baseline F1-score by 9% on cross-validation. They experimented with stylistic features, LiX Score, TF-IDF Unigrams, TF-IDF Profanity, TF-IDF Emojis, POS Tag Counts, Sentiment Analysis based features, count of punctuation marks[29]. Of the classification algorithms considered, they find that the Random Forest classifier performs the best, achieving a final F1-score of 96.04% with a 70/30 split on the train set and a final accuracy of 95.56% on the test data provided by PAN.

Marco Siino et al., proposed [15] a novel ensemble model based on deep learning and non-deep learning classifiers. The proposed model was developed for participating at the Profiling Irony and Stereotype Spreaders (ISSs) task hosted at PAN@CLEF2022. Our ensemble (named T100), include a Logistic Regressor (LR) that classifies an author as ISS or not (nISS) considering the predictions provided by a first stage of classifiers. All these classifiers are able to reach state-of-the-art results on several text classification tasks. These classifiers are a Convolutional Neural Network (CNN), a Support Vector Machine (SVM), a Decision Tree (DT) and a Naive Bayes (NB) classifier. The classifiers or voters are trained on the provided dataset and then generate predictions on the training set. Finally, the LR is trained on the predictions made by the voters. For the simulation phase, the LR considers the predictions of the voters on the unlabelled test set to provide its final prediction on each sample. To develop and test their model, they used a 5-fold cross validation on the labelled training set. Over the five validation splits, the proposed model achieves a maximum accuracy of 0.9342 and an average accuracy of 0.9158. As announced by the task organizers, the trained model presented here is able to reach an accuracy of 0.9444 on the unlabelled test set provided for the task.

The use of stereotypes, irony, mocking and scornful language is prevalent on social media platforms such as Twitter. Identification or profiling of users who are involved in the spread of such content is beneficial for monitoring its spread. Dhaval Taunk et al., studied [16] the problem of profiling irony and stereotype spreaders on Twitter as a part of the PAN shared task in CLEF 2022. Their experimentation pipeline consisting of pre-processing followed by featured extraction and ML-based modelling. A pre-processed tweet text is represented by using a TF-IDF vector as a simple term-frequency based tweet representation and the representation of a user is calculated as a summation of all the tweet vectors. They experimented with classifiers based on Logistic Regression, K Nearest Neighbours, Support Vector Machines, Random Forest and XGBoost. The best performing Random Forest Classifier was used to take predictions on the test set and was submitted for evaluation on the TIRA platform, which gave an accuracy of 95%.

The study of irony detection on social networks has gained much attention in recent years. As part of this task, a collection of users' tweets is provided, and the goal is to determine if these users are spreaders of irony or not. In [17], as they hypothesized that user-generated content is affected by the author's psychometric characteristics, contextual information, and irony features in the text. Authors investigated the effects of incorporating this information to identify ironic spreaders. Using the emotion and personality detection module, they were able to extract the author's psychometric features. A pre-trained language model based on SBERT and T5-based architecture has been employed to extract context-based features. The proposed work describes a framework by using the author's psychometric, contextual, and ironic features in a Gradient Boosting classifier. Experimental results of their work demonstrate the importance of this combination in identifying ironic spreader users. As a result, they were able to achieve an accuracy of 95.00% and 93.81% with 5-fold and 10-fold cross-validation respectively on the IROSTEREO training dataset. However, on the official PAN test set, our system attained an 88.89% score.

III. DATASET CHARACTERISTICS

PAN organizing competitions on different task every year by providing datasets in different languages [18]. Participants from around the globe are submitting their works to the competition. Author Profiling is one task started in the year 2013 to predict the demographic characteristics like age, gender, nativity language and personality traits in successive years [19]. In 2022 PAN competition, the organizers introduced a variety of author profiling task of irony and stereotypes spreaders detection. Unlike the previous year, this time only texts in English were used, with a selection of 420 authors for training, distributed into 210 ironic and 210 non-ironic profiles, which represents a balanced data corpus. The dataset is modelled in such a way that each user can be profiled as Irony spreader (I) and Not Irony spreader (NI). The organizers also provided the participant with 180 blinded XML files (each contains 200 tweets) as test set for profiling Irony spreaders. The Table 1 shows the characteristics of dataset.

Table 1: Statistics about the training data

	Irony(I)	Not Irony (NI)
No. of Authors	210	210
No. of Tweets	42000	42000
Mean Length	5979	5296
Length of Vocabulary	57247	46642

IV. MACHINE LEARNING ALGORITHMS

In this work, we used two machine learning algorithms such as Support Vector Machine (SVM) and Random Forest to train the document vectors and for predicting the accuracy of proposed method.

4.1 Support Vector Machine

In 1963, SVMs were originally proposed by Vapnik to solve binary classification problems. Recently, SVMs have been extended to adapt to multi-class problems [20]. The main goal of SVMs is to build or construct one or more hyperplanes to split a given dataset into multiple subsets corresponding to different class labels. There might be many hyperplanes that can split the data but the selected hyperplane should maximize the distances with its nearest training instances.

In many practical problems, SVMs achieve good performance [21]. However, users usually have to provide a good kernel function for SVMs for non-linear cases. SVM is able to handle sparsity of the data representation by using right Kernel functions or regularization methods. The most popular kernels available are linear, polynomial, RBF (Radial Basis Function) and sigmoid. In terms of efficiency, SVMs have a high computational cost and require a large memory in the training phase when there is more number of dimensions [22].

4.2 Random Forest

Random Forest consists of a large number of individual decision trees that operate as an ensemble. Each individual tree in the random forest spits out a class prediction and the class with the most votes becomes our model prediction [23]. This model performs well for given large feature sets because it combines the predictions of various decision trees to build a more robust classifier. While constructing new decision trees, this method uses a random subset of features which gets rid of spurious features and improving the robustness of our estimate. The random forest method operated by following a sequence of steps such as selection of random samples from the given dataset, generate a decision tree for each sample, gather prediction result of each decision tree, perform voting based on all outcomes of decision trees, and predict final class based on the majority votes.

In random forest algorithm, different hyperparameters are used to either improve the model's performance and predictive capacity or to make it faster. The hyperparameters are estimators (the number of trees built by the algorithm before averaging the predictions), Max_features (the number of features that a random forest evaluates while splitting a node), mini_sample_leaf (the number of leaves are necessary to separate an internal node), n_jobs (number of processors it can use), random state (regulates the sample's unpredictability), and oob score (It's a cross-validation method based on random forests).

V. EVALUATION MEASURES

The evaluation measures are used by the machine learning algorithms to estimate the efficiency of the proposed system. The researchers used various measures like recall, precision, F1-score as well as the accuracy to check the efficiency of the developed system [24]. The Precision is the percentage of documents that the classifier labels as relevant that is actually relevant. Equation (1), (2), (3), and (4) are used to calculate Precision, recall, F1-Score, and accuracy respectively.

$$\text{Precision} = \frac{TP}{TP + FP} \quad (1)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (2)$$

$$\text{F1-Score} = 2 \times \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3)$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (4)$$

Where, TP is number of given documents classifies as positive and is also in the actual positive class, FP is number of given documents classifies as positive but is in the actual negative class, FN is number of given documents classifies as negative but is in the actual positive class, TN is number of given documents classifies as negative and is also in the actual negative class.

The accuracy measure is used in this work to present the efficiency of our proposed approach.

6. PROPOSED NOVEL TEXT REPRESENTATION METHOD

In this work, we proposed a novel text representation method for irony and stereotypes spreaders detection in Twitter. The model for the proposed method is represented in Figure 1.

In the proposed method, the first and foremost important step is collection of dataset for experimentation. In this work, we gathered the dataset for irony and stereotypes spreaders detection from the PAN 2022 competition. Once the dataset is collected the next step is prepare the dataset for extracting required features from the dataset. Pre-processing techniques are applied to prepare the dataset for feature extraction. In this work, we applied different pre-processing techniques such as remove the generic hashtag, user and URL tags from the tweets, Converts all text to lowercase, Convert Emoji's to the corresponding text, Convert sentence abbreviations to extended mode, Delete Duplicate words and simplify. Then tokenize the data by using the TweetTokenizer provided by NLTK. After cleaning the irrelevant data from the dataset, next step is extracting all informative terms from the dataset. All the terms are represented as vectors. These term vectors are used to represent the documents as vectors. The document vectors are trained with machine learning algorithms to generate the accuracy of proposed method for irony and stereotypes spreaders detection.

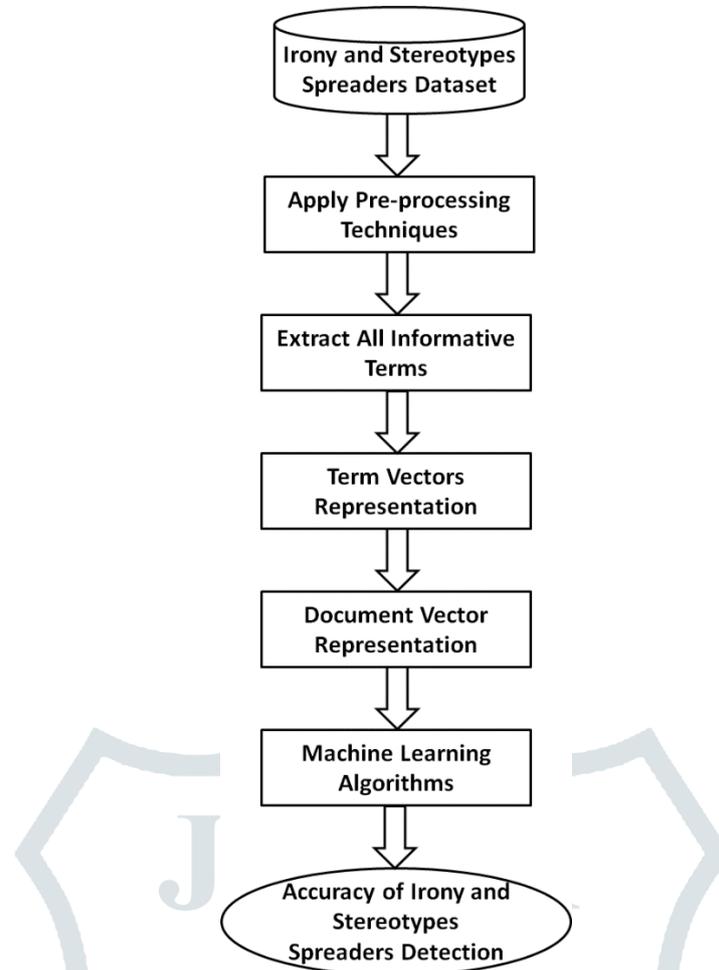


Figure 1: The Model for novel text representation method

6.1 Term Vector Representation

All the terms those are extracted from the dataset are represented as vectors. Table 2 shows the representations of term vectors.

Table 2: Representations of Term Vectors

	D_1	D_2	D_m
T_1	W_{11}	W_{12}	W_{1m}
T_2	W_{21}	W_{22}	W_{2m}
...
...
T_n	W_{n1}	W_{n2}	W_{nm}

In Table 2, $\{T_1, T_2, \dots, T_n\}$ is the set of terms extracted from the dataset. $\{D_1, D_2, \dots, D_m\}$ is the set of documents in the dataset. W_{nm} denotes the weight of term T_n in document D_m . The weight of term in a document is computed by using the term weight measure. In this work, we proposed a combined term weight measure of TF-IDF-RF-ICSDF for computing term weight. Each term is represented as an m-dimensional vector.

6.2 Document Vector Representation

The document is represented as a vector by using the terms in that document. Each document is represented as an m-dimensional vector. The vector of document is generated by aggregating the term vectors that are contained in that document. The Table 3 shows the vector representations of documents.

Table 3: Representations of document vectors

Documents	Terms in Document	Term Vector				
		1	2	m
D ₁	T ₁	W ₁₁	W ₁₂	W _{1m}
	T ₂	W ₁₂	W ₂₂	W _{2m}

	T _i	W _{i1}	W _{i2}	W _{im}
D ₁ Vector Representation		Average of all Dimension-1	Average of all Dimension-2	Average of all Dimension-m
D ₂	T ₁	W ₁₁	W ₁₂	W _{1m}
	T ₂	W ₁₂	W ₂₂	W _{2m}

	T _j	W _{j1}	W _{j2}	W _{jm}
D ₂ Vector Representation		Average of all Dimension-1	Average of all Dimension-2	Average of all Dimension-m
...
...
D _m	W ₁₁	W ₁₂	W _{1m}	W ₁₁
	W ₁₂	W ₂₂	W _{2m}	W ₁₂

	W _{i1}	W _{k2}	W _{km}	W _{k1}
D _m Vector Representation		Average of all Dimension-1	Average of all Dimension-2	Average of all Dimension-m

In Table 3, the document D₁ contains ‘i’ number of terms, document D₂ contains ‘j’ number of terms, and document D_m contains ‘k’ number of terms. The document D₁ is represented as a m-dimensional vector by aggregating the ‘i’ number of m-dimensional vectors. Likewise, document D_m is represented as a m-dimensional vector by aggregating the ‘k’ number of m-dimensional vectors.

6.3 Term Weight Measure (TF-IDF-RF-ICSDF)

In this work, we proposed a term weight measure TF-IDF-RF-ICSDF by combining different factors of different term weight measures. The TF measure assigns more weight to terms that are occurred more times in a document [25][28]. IDF measure assign more weight to the terms are occurred in less number of documents in whole dataset [25]. RF measure gives more weight to the terms that are occurred in more positive class documents than negative class documents. ICSDF measure allots more weight to the terms that are occurred in less number of documents and that are distributed in less number of classes [26]. The following Equation is used to determine the weight of a term T_i in a document D_k by using TF-IDF-RF-ICSDF measure.

$$TF - IDF - RF - ICSDF(T_i, D_k) = TF(T_i, D_k) \times \left(1 + \log\left(\frac{N}{DF(T_i) + 1}\right)\right) \times \log\left(2 + \frac{A}{C}\right) \times \left(1 + \log\left(\frac{m}{\sum_{j=1}^m \left(\frac{n_{cj}(T_i)}{N_{cj}}\right)}\right)\right)$$

Where, N is count of documents in dataset, $DF(T_i)$ is count of documents contain term T_i in total dataset, m is count of classes in the dataset, $n_{c_j}(T_i)$ is count of documents in class j^{th} class contain term T_i , N_{c_j} is count of documents in j^{th} class. A is Number of positive class documents contain term T_i , C is Number of negative class documents contain term T_i .

7. EXPERIMENTAL RESULTS

In this work, the experiment performed for predicting the accuracy of irony and stereotypes spreaders detection. The proposed text representation method converts each document as a vector of m -dimensional. These document vectors are trained with two machine learning algorithms such as support vector machines and random forest. The performance of two machine learning algorithms for predicting the accuracy of irony and stereotypes spreaders detection is depicted in Table 4.

Table 4: Accuracies of Proposed Method for irony and stereotypes spreaders detection

Machine Learning Algorithms	Accuracy of irony and stereotypes spreaders detection
Support Vector Machine	94.28
Random Forest	97.86

In Table 4, the proposed text representation method attained best accuracies of 94.28 and 97.86 for irony and stereotypes spreaders detection when experimented with the machine learning algorithms of support vector machine and random forest respectively. We conducted experiment with other machine learning algorithms such as K-Nearest Neighbour (KNN), Naive Bayes Multinomial (NBM) etc., but we observed that the SVM and RF algorithms gave good accuracies for irony and stereotypes spreaders detection. We also observed that the random forest algorithms shows best performance for irony and stereotypes spreaders detection when compared with other machine learning algorithms.

8. CONCLUSIONS AND FUTURE SCOPE

In this work, we proposed a novel text representation method for irony and stereotype spreaders detection. In the proposed method extract all informative terms from the dataset and represent each term as a vector by computing weight of a term specific to database documents. Documents are represented as vector by using the terms that are contained in that document. Two machine learning algorithms such as SVM and RF are used for presenting the performance of proposed method. The proposed text representation method attained best accuracies of 94.28 and 97.86 for irony and stereotypes spreaders detection when experimented with the machine learning algorithms of support vector machine and random forest respectively.

In future works, we would like to examine different language models over the proposed framework and other architectures to achieve the best performance of the current framework. Also, we plan to investigate a variety of features like sentiment analysis to boost the current text representation method.

REFERENCES

- [1] Raghunadha Reddy T, Vishnu Vardhan B, Vijayapal Reddy P, "A Survey on Author Profiling Techniques", International Journal of Applied Engineering Research, March 2016, Volume-11, Issue-5, pp. 3092-3102.
- [2] Swathi Ch, Karunakar K, Archana G, T. Raghunadha Reddy, "A New Term Weight Measure for Gender Prediction in Author Profiling", Proceedings in Advances in Intelligent Systems and Computing, Volume 695, PP. 11-18, 2018.
- [3] Francisco Rangel and Paolo Rosso. Overview of the 7th Author Profiling Task at PAN2019: Bots and Gender Profiling. In Linda Cappellato, Nicola Ferro, David E. Losada, and Henning Müller, editors, CLEF 2019 Labs and Workshops, Notebook Papers, September 2019. CEUR-WS.org.
- [4] Francisco Rangel, Anastasia Giachanou, Bilal Ghanem, and Paolo Rosso. Overview of the 8th Author Profiling Task at PAN 2020: Profiling Fake News Spreaders on Twitter. In Linda Cappellato, Carsten Eickhoff, Nicola Ferro, and Aurélie Névéol, editors, CLEF 2020 Labs and Workshops, Notebook Papers, September 2020. CEUR-WS.org.
- [5] Francisco Rangel, Paolo Rosso, Gretel Liz De La Peña Sarracén, Elisabetta Fersini, and BERTa Chulvi. Profiling Hate Speech Spreaders on Twitter Task at PAN 2021. In Guglielmo Faggioli, Nicola Ferro, Alexis Joly, Maria Maistro, Florina Piroi, editors, CLEF 2021 Labs and Workshops, Notebook Papers, 2021. CEUR-WS.org.
- [6] Janek Bevendorff and Berta Chulvi and Elisabetta Fersini and Annina Heini and Mike Kestemont and Krzysztof Kredens and Maximilian Mayerl and Reyner Ortega-Bueno and Piotr Pezik and Martin Potthast and Francisco Rangel and Paolo Rosso and Efstathios Stamatatos and Benno Stein and Matti Wiegmann and Magdalena Wolska and Eva Zangerle: Overview of PAN 2022: Authorship Verification, Profiling Irony and Stereotype Spreaders, and Style Change Detection. In: Proceedings of the Thirteenth International Conference of the CLEF Association (CLEF 2022)
- [7] Dr. T. Raghunadha Reddy, M. Shashi Preetham, K. Sree Vasini, A. Rajesh, "A Deep Learning Approach for Hate Speech Spreaders Detection using Statistical and Contextualized Embeddings", Journal of Emerging Technologies and Innovative Research (JETIR), Volume 10, Issue 5, May 2023.
- [8] Maria Fernanda Artigas-Herold, Daniel Castro-Castro, "User profiling: voting scheme", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy

- [9] Sabur Butt, Fazlourrahman Balouchzahi, Grigori Sidorov and Alexander Gelbukh, "CIC@PAN: Simplifying Irony Profiling using TwitterData", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [10] Daniele Croce, Domenico Garlisi and Marco Siino, "An SVM Ensemble Approach to Detect Irony and Stereotype Spreaders on Twitter", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [11] José Antonio García-Díaz, Miguel Ángel Rodríguez-García, Francisco García-Sánchez and Rafael Valencia-García, "UMU Team at IROSTEREO: Profiling Irony and Stereotype spreaders on Twitter combining Linguistic Features with Transformers", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [12] Catherine Ikae, "UniNE at PAN-CLEF 2022: Profiling Irony and Stereotype Spreaders on Twitter", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [13] Hyewon Jang, "Lexicon-Based Profiling of Irony and Stereotype Spreaders", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [14] Tiago Filipe Nunes Ribeiro, Yana Nikolaeva Nikolova and Kaja Seraphina Elisa Hano, "Irony & Stereotype Spreader Detection using Random Forests", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [15] Marco Siino, Ilenia Tinnirello and Marco La Cascia, "T100: A modern classic ensemble to profile irony and stereotype spreaders", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [16] Dhaval Taunk, Sagar Joshi and Vasudeva Varma, "Profiling irony and stereotype spreaders on Twitter based on term frequency in tweets", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [17] Ehsan Tavan*1, Maryam Najafi*1 and Reza Moradi, "Identifying Ironic Content Spreaders on Twitter using Psychometrics, Contextual and Ironic Features with Gradient Boosting Classifier", Notebook for PAN at CLEF 2022, CLEF 2022 – Conference and Labs of the Evaluation Forum, September 5–9, 2022, Bologna, Italy.
- [18] Raghunadha Reddy T, Vishnu Vardhan B, Vijayapal Reddy P, "Profile specific Document Weighted approach using a New Term Weighting Measure for Author Profiling", *International Journal of Intelligent Engineering and Systems*, 9 (4), pp. 136-146, Nov 2016.
- [19] Raghunadha Reddy T, Vishnu Vardhan B, Vijayapal Reddy P, "Author profile prediction using pivoted unique term normalization", *Indian Journal of Science and Technology*, Vol 9, Issue 46, Dec 2016.
- [20] Cortes, C. & Vapnik, V. *Machine Learning* (1995) 20: 273. <https://doi.org/10.1023/A:1022627411411>
- [21] Raghunadha Reddy T, Vishnu Vardhan B, Vijayapal Reddy P, "A Document Weighted Approach for Gender and Age Prediction", *International Journal of Engineering -Transactions B: Applications*, Volume 30, Number 5, pp. 647-653, May 2017.
- [22] Archana Gelli, Karunakar Kavuri, T Raghunadha Reddy, Lakshmi Narayana M, "Distance Measures based Approach for Hate Speech Spreaders Detection", *Journal of Applied Science and Computations*, Volume IX, Issue XII, December/2022, Pages: 227 – 233
- [23] Raghunadha Reddy T, Vishnu Vardhan B, GopiChand M, Karunakar K, "Gender prediction in Author Profiling using ReliefF Feature Selection Algorithm", *Proceedings in Advances in Intelligent Systems and Computing*, Volume 695, PP. 169-176, 2018.
- [24] T. Raghunadha Reddy, P. Vijayapal Reddy, T Murali Mohan, Raju Dara, "An Approach for Suggestion Mining based on Deep Learning Techniques", *International Conference on Computer Vision, High Performance Computing, Smart Devices and Networks (CHSN-2020)*, 28-29 December, 2020, JNTUK, Kakinada, Andhra Pradesh. IOP Conference Series: Materials Science and Engineering, DOI 10.1088/1757-899X/1074/1/012021, Volume 1074
- [25] T. Raghunadha Reddy, P. Vijaya Pal Reddy, P. Chandra Sekhar Reddy, "A New Supervised Term Weight Measure based Machine Learning Approach for Text Classification", *International Conference on Intelligent Systems and Sustainable Computing*, September 24-25, 2021, Malla Reddy University Hyderabad, India, *Intelligent Systems and Sustainable Computing, Smart Innovation, Systems and Technologies*, pp 563–571, ISBN : 978-981-19-0011-2, vol 289
- [26] T. Raghunadha Reddy, "A Novel Approach for Authorship Verification using Similarity Measure", *International Journal for Innovative Engineering and Management Research*, Volume 09, Issue 12, 2020, pp. 437-445.
- [27] Kavuri, K. ., & Kavitha, M. (2023). A Word Embeddings based Approach for Author Profiling: Gender and Age Prediction. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11(7s), 239–250. <https://doi.org/10.17762/ijrctc.v11i7s.6996>.
- [28] K. Kavuri and M. Kavitha, "A Term Weight Measure based Approach for Author Profiling," *2022 International Conference on Electronic Systems and Intelligent Computing (ICESIC)*, Chennai, India, 2022, pp. 275-280, doi: 10.1109/ICESIC53714.2022.9783526.
- [29] Karunakar Kavuri, Kavitha, M. (2020). "A Stylistic Features Based Approach for Author Profiling". In: Sharma, H., Pundir, A., Yadav, N., Sharma, A., Das, S. (eds) *Recent Trends in Communication and Intelligent Systems. Algorithms for Intelligent Systems*. Springer, Singapore. https://doi.org/10.1007/978-981-15-0426-6_20.