# Cyberspace News Prediction of Text and Image with Report Generation

**Tallapudi Jyothsna**

**PG** scholar, Department of MCA, DNR College, Bhimavaram, Andhra Pradesh.

**K.Sridevi**

(Assistant Professor), Master of Computer Applications, DNR college, Bhimavaram, Andhra Pradesh.

*Abstract: The cyberspace news consumption is increasing day by day all over the world. The main reason for cyber space news consumption is due to its rapid spread of information and its easy access which lead people to consume news rapidly without the knowledge of whether the news is false or true. Thus, it leads to the wide spread of false news which leads to the negative impacts on society. Therefore false news prediction on cyberspace is attracting a tremendous attention. The issue of fake-news prediction on cyberspace is both challenging and relevant as spreading of fake news occurs in various streams like text, audio, video, images etc. This model works on processing the text and images together by providing an interactive Application Interface (API), i.e. text by applying the model Logistic regression classifier and image by applying self-consistency algorithm. The natural language tool kit (NLTK) model is used for these implementation through python. Once the news is predicted fake, a report is redirected to the authorized website (cybercrime department) to take the immediate necessary actions required to stop these news from spreading.*

*Index Terms—Cyberspace, fake-news, text and image, Logistic regression classifier, self-consistency algorithm, report, redirect*

## I. INTRODUCTION

The aim of a study on "Cyberspace News Prediction of Text and Image with Report Generation" is to develop a comprehensive and advanced system that can predict and generate news reports in the cyberspace domain by utilizing both text and image data

### Objective

The objectives of a study on "Cyberspace News Prediction of Text and Image with Report Generation" are specific goals and tasks that need to be achieved in order to fulfill the broader aim of the research. These objectives provide a roadmap for conducting the study and building the proposed system effectively.

### Scope

The scope of a study on "Cyberspace News Prediction of Text and Image with Report Generation" encompasses the boundaries, limitations, and the range of activities and aspects that the research project will cover.

### Introduction

**N**OWADAYS, people spend a lot of time in Internet (cyberspace) and consume news. The main reason for rapid spread of news in cyberspace is due to its low cost, easy access and easy sharing facility. This made people to consume news from cyberspace rather than fetching it from television or newspaper. The widespread of fake-news will have a serious negative impact on society and individuals.

Fake-news detection on cyberspace has led to tremendous research all over the world to predict with the exact accuracy as the content of false-news is diverse in topics. People consuming news from cyberspace produce data which is diverse and difficult to predict This model is a solution to all these problems of fake news in cyberspaces that is fast growing. In particular the datasets which are trained by various machine learning techniques like data pre-processing, feature selection, self-consistency etc. and all these are implemented by natural language processing in python.

Here we detect both forms of fake news, i.e., both text and image streams. Once the prediction is false the report is generated and it is immediately redirected to the authorized page (cybercrime department) insisting the seriousness of the news for which the actions will be taken accordingly. Through this we try to bring a safe and trustable cyberspace experience to people who rely on this. They can now verify news before they are believing or forwarding them to others.

### CHAPTER 2 LITERATURE SURVEY

**Tallapudi Jyothsna/ International Journal of Management Research & Review**

Social networking sites engage millions of users around the world. The users' interactions with these social sites, such as Twitter and Facebook have a tremendous impact and occasionally undesirable repercussions for daily life. The prominent social networking sites have turned into a target platform for the spammers to disperse a huge amount of irrelevant and deleterious information. Twitter, for example, has become one of the most extravagantly used platforms of all times and therefore allows an unreasonable amount of spam.

Fake users send undesired tweets to users to promote services or websites that not only affect legitimate users but also disrupt resource consumption. Moreover, the possibility of expanding invalid information to users through fake identities has increased that results in the unrolling of harmful content. Recently, the detection of spammers and identification of fake users on Twitter has become a common area of research in contemporary online social Networks (OSNs). In this paper, we perform a review of techniques used for detecting spammers on Twitter.[1]

Twitter is rated as the most popular social network among teenagers [2]. However, exponential growth of Twitter also invites more unsolicited activities on this platform. Nowadays, 200 million users generate 400 million new tweets per day [3]. This rapid expansion of Twitter platform influences more number of spammers to generate spam tweets which contain malicious links that direct a user to external sites containing malware downloads, phishing, drug sales, or scams [4]. These types of attacks not only interfere with the user experience but also damage the whole internet which may also possibly cause temporary shutdown of internet services all over the world

An accurate and efficient face recognition system is a more interesting topic in most industries and research areas. It is a type of biometric information process that is easily adaptable as compared to the tradition card recognition system. Generally, a face recognition system is preceded by a face detection technique. The face detection technique is the preliminary stage to detect a face in live images.

In this paper, some face detection techniques are discussed such as finding skin likelihood image, skin segmentation, the morphological operation for extracting boundary regions, Haarlike features, and Ada-boost algorithm. This Haar-like feature algorithm continually searches its pattern from the particular face and which has better advantages over other techniques. After the face detection technique, the face recognition technology is applied on the detected face for further identification by using some classifiers.

**CHAPTER 4**

**PROPOSED METHOD**

Upload news dataset

↓

Preprocess dataset

↓

TF-IDF Vector generation

↓

Run NB algorithm

↓

Run SVM algorithm
Run LR algorithm

↓

Run RF algorithm
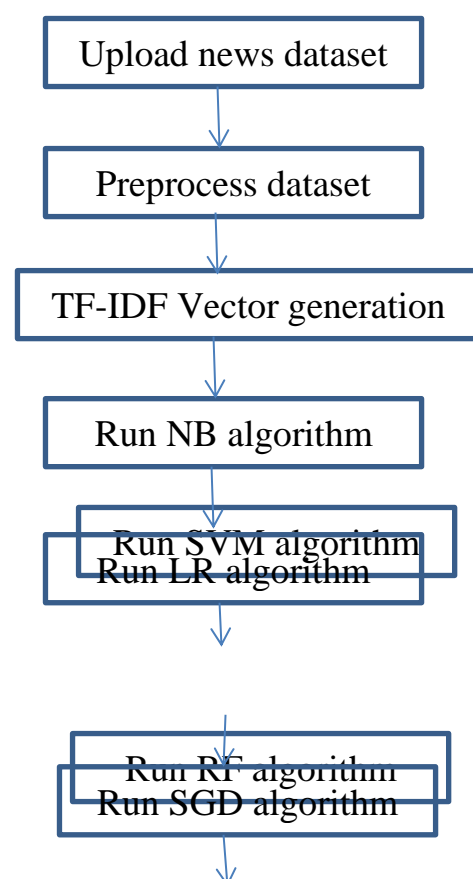Run SGD algorithm

↓

**Fig. 1. Flowchart**

To implement this project we have designed following modules.

1) Upload News Dataset: using this module we will upload dataset directory with REAL and FAKE news files to application

2) Preprocess Dataset: using this module we will read all text news and then apply NLTK technique to clean text news data

3) TF-IDF Vector Generation: cleaned text news will be input to TF-IDF algorithm to convert text data into numeric vector and then convert vector into train and test where application used 80% dataset for training and 20% for testing

4) Run Naive Bayes Algorithm: 80% train data will be input to Naïve Bayes algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy. The higher the accuracy the better is the algorithm

5) Run Logistic Regression Algorithm: 80% train data will be input to Logistic Regression algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

6) Run SVM Algorithm: 80% train data will be input to SVM algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

7) Run SGD Algorithm: 80% train data will be input to SGD algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

8) Run Random Forest Algorithm: 80% train data will be input to Random Forest algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

9) Comparison Graph: using this module we will plot accuracy graph of all algorithms

10) Predict Fake Text News: using this module user can enter his text and then ML algorithm will predict weather given news in TRUE or FALSE

11) Predict Fake Image News: using this module we will upload IMAGES and then visual analysis algorithm will predict weather image is TRUE or FALSE

## CHAPTER 5

## RESULT

Cyberspace News Prediction of Text and Image with Report Generation

In this paper author is using various machine learning algorithms such as Naïve Bayes, Logistic Regression, SVM, SGD and Random Forest for fake news prediction on TEXT data and then applying VISUAL content analysis algorithm on images to predict it as fake or real news images. Most of the time fake news images are the duplicate copy of old original images so by analysing visual content of old and new images we can predict weather image uses in NEWS is real or fake.

To analyse text news author is using NLTK technique to remove stop words, special symbols, and applying stemming and lemmatization to clean text news and then converting clean text news into features vector by applying TF-IDF algorithm. TF-IDF will replace each word with its average frequency to build TF-IDF vector. TF-IDF extracted features will be input to all ML algorithms to train a model and in all algorithms Random Forest and Logistic Regression is giving better accuracy.

In this paper author is using fake and real news dataset with label as TRUE and False. Both TF-IDF features and labels will be input to all algorithms to build a model and this model will be applied on TEXT data to predict image as REAL or FAKE.

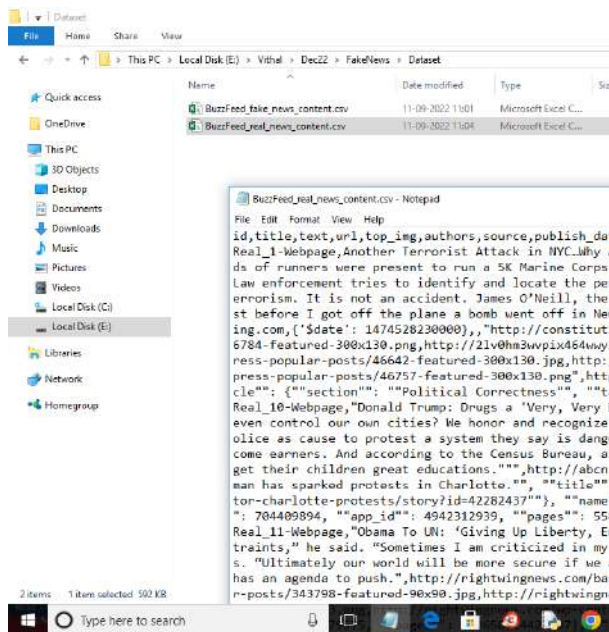We are using below dataset for this project

**Fig.2 Dataset**

In above dataset we have two folders where one will contains REAL news and other contains FAKE news and from above dataset we are extracting 'TEXT' column as NEWS.

To implement this project we have designed following modules.

1. Upload News Dataset: using this module we will upload dataset directory with REAL and FAKE news files to application

2. Preprocess Dataset: using this module we will read all text news and then apply NLTK technique to clean text news data

3. TF-IDF Vector Generation: cleaned text news will be input to TF-IDF algorithm to convert text data into numeric vector and then convert vector into train and test where application used 80% dataset for training and 20% for testing

4. Run Naive Bayes Algorithm: 80% train data will be input to Naïve Bayes algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy. The higher the accuracy the better is the algorithm

5. Run Logistic Regression Algorithm: 80% train data will be input to Logistic Regression algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

6. Run SVM Algorithm: 80% train data will be input to SVM algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

7. Run SGD Algorithm: 80% train data will be input to SGD algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

8. Run Random Forest Algorithm: 80% train data will be input to Random Forest algorithm to train a model and this model will be applied on 20% test data to calculate prediction accuracy

9. Comparison Graph: using this module we will plot accuracy graph of all algorithms

10. Predict Fake Text News: using this module user can enter his text and then ML algorithm will predict weather given news in TRUE or FALSE

11. Predict Fake Image News: using this module we will upload IMAGES and then visual analysis algorithm will predict weather image is TRUE or FALSE

To run project double click on 'run.bat' file to get below screen

**Fig.3 Upload News Dataset**

In above screen click on 'Upload News Dataset' button to upload NEWS dataset and get below output
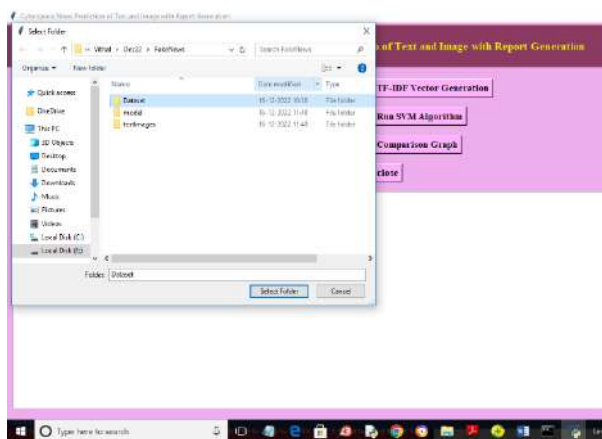


**Fig.4 selecting and uploading entire 'Dataset' folder**

In above screen selecting and uploading entire 'Dataset' folder and then click on 'Select Folder' button to load dataset and get below output
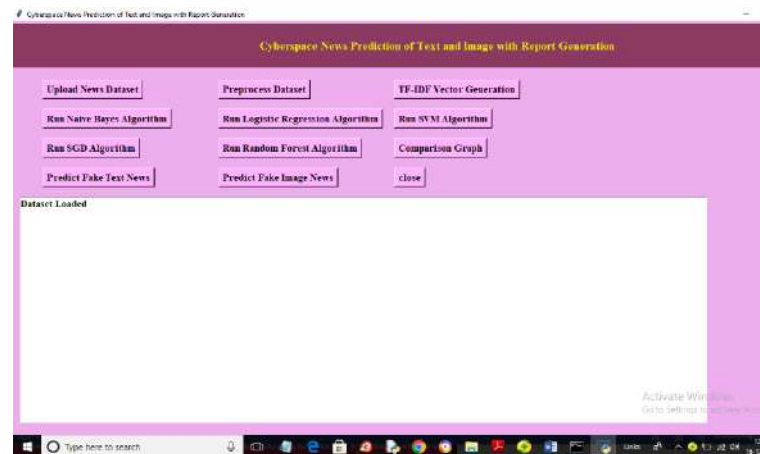


**Fig.5 dataset loaded**

In above screen dataset loaded and now click on "Preprocess Dataset' button to read and clean dataset and get below output



**Fig.6 clean news text loaded**

In above screen clean news text loaded and in above text stop words, special symbols are removed and then applied stemming and lemmatization technique. Now click on 'TF-ID Vector Generation' button to generate vector and get below output
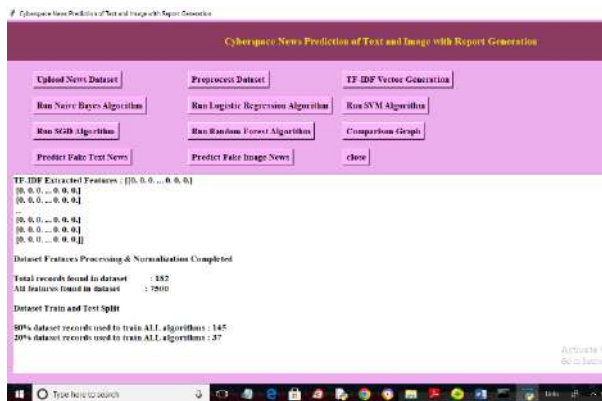
**Fig.7 TF-IDF vector generated**

In above screen TF-IDF vector generated and then we can see dataset contains 182 news and each news contains 7500 features and then we split news into train and test and now click on 'Run Naïve Bayes Algorithm' button to train Niave Bayes and get below accuracy



**Fig.8 Naïve Bayes we got 97% accuracy**

In above screen with Naïve Bayes we got 97% accuracy and similarly click all algorithms button to get below output

**Fig.9. all algorithms got accuracy more than 95%**

In above screen we can see all algorithms got accuracy more than 95% and this accuracy may vary for each run as we are splitting train and test data randomly. Now click on 'Comparison Graph' button to get below output
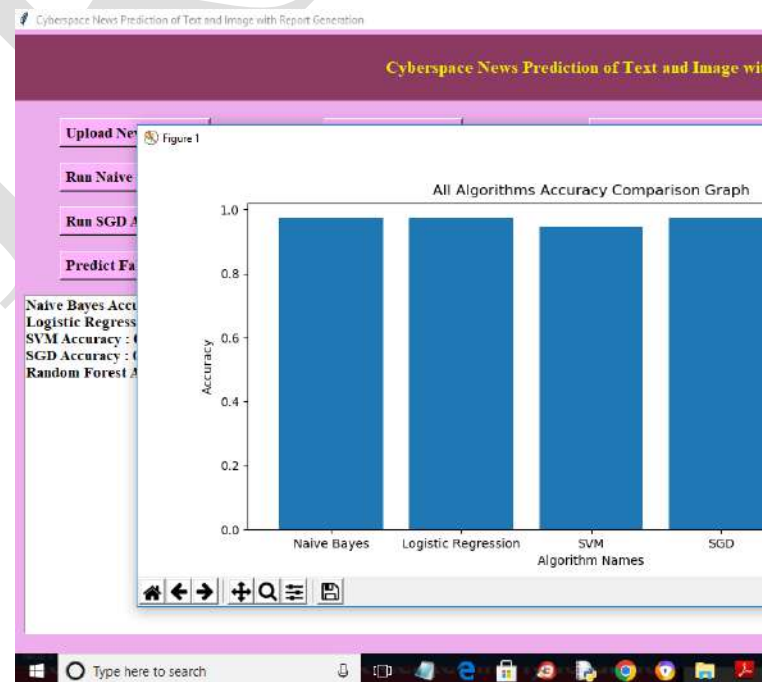


**Fig.10 Graphical representation**

In above graph x-axis represents algorithm names and y-axis represents accuracy of the algorithms. Now close above graph and then click on 'Predict Fake Text News' button to enter text news and get prediction output
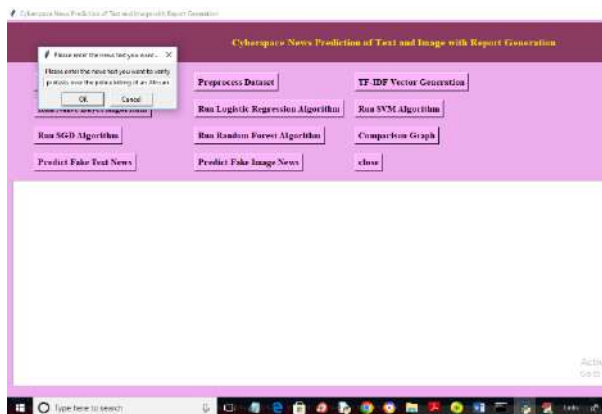
205

**Fig.11 dialog box I entered some news**

In above screen in dialog box I entered some news and press 'OK' button to get below output



**Fig.12 displaying entered news**

In above screen I am displaying entered news and then displaying prediction output as news is 'False or Fake' and truth probability is 0.24%. Now I will enter some other news to get below output
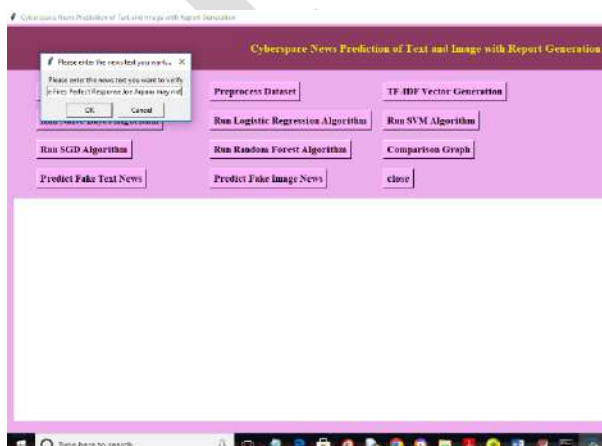


**Fig.13 entered some other news**

In above screen I entered some other news and press OK button to get below output



**Fig.14 entered news predicted as 'TRUE'**

In above screen entered news predicted as 'TRUE' and similarly you can enter some news and get prediction. If you want you can copy and paste news from 'Test_News.csv' file. Now click on 'Predict Fake Image News' button to upload image and get below output
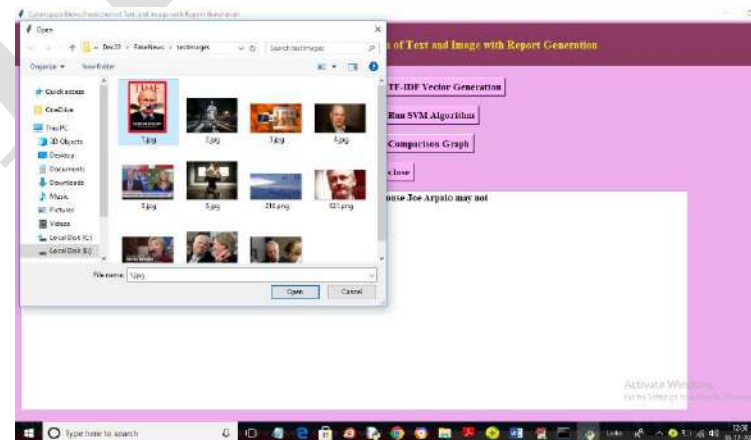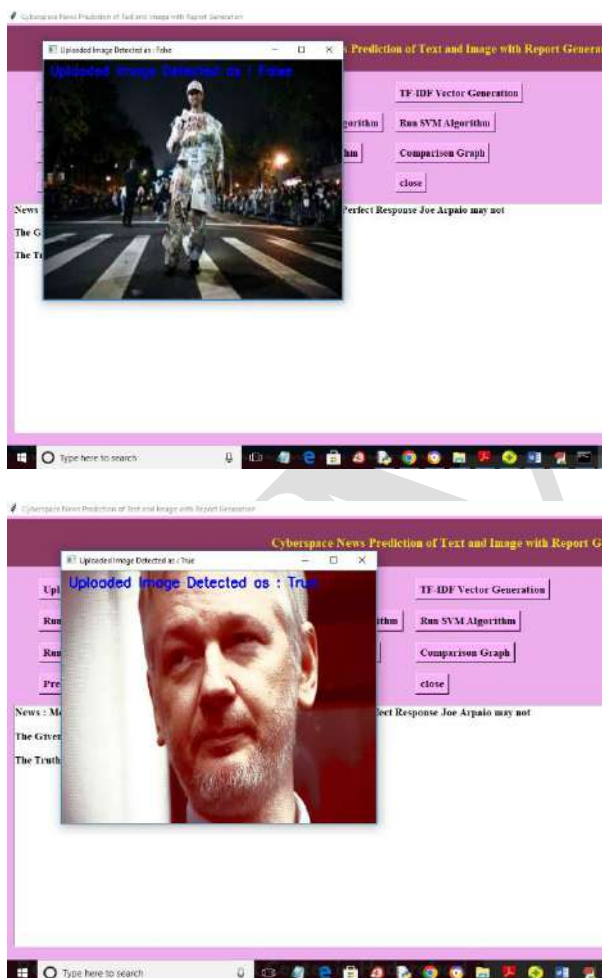


**Fig.15 uploading image**

In above screen uploading image and then press 'Open' button to to get below output

**Fig.16 image predicted as False News**

In above screen image predicted as False News and similarly you can upload and test other images





**CHAPTER 6**

**CONCLUSION**

The consumption of news is increasing day by day in cyberspace than the traditional media. Due to its increasing popularity and user friendly access it leaves a huge impact on individuals and society. Therefore, in this model we have found a way to detect such fake news in both the forms of text and image by using the Logistic regression model. By redirecting the fake news to the authorized website (cybercrime department), we hereby frame a high social impact and thus it reduces the spreading of false news distinctly. This model can be further discussed for the future improvement in fake news detection which can be in audio, video streams and commercialize the field to other applications.

**REFERENCES**

[1] Faiza Masood, Ghana Ammad, Ahmad Almogren, Assad Abbas, Hasan Ali Khattak, Ikram Ud Din, Mohsen Guizani and Mansour Zuair, "Spammer Detection and Fake User Identification on Social Networks," IEEE Trans. Inf. Translations and content mining, vol. 7, pp. 2169- 3536, 2019.

[2] Himank Gupta, Mohd. Saalim Jamal, Sreekanth Madisetty and Maunendra Sankar Desarkar, "A framework for realtime spam detection in Twitter," IEEE Int. Conf. Communication Systems and networks, pp. 2155-2509, 2018.

[3] K.Sakthidasan, G.Srinithya, V.Nagarajan (FEB 2014), "Enhanced Edge Preserving Restoration for 3D Images Using Histogram Equalization Technique", International Journal of Electronic Communications Engineering Advanced Research, Vol.2, SP-1, Feb.2014, pp. 40-44

[4] S. Kwon, M. Cha, K. Jung, W. Chen and Y. Wang, "Prominent features of rumor propagation in online social media," IEEE Int. Conf. Data Mining, pp. 1103–1108, 2013.

[5] Hadeer Ahmed, Issa Traore and Sherif Saad, "Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques," Springer, pp. 127–138, 2017.

[6] K. Wu, S. Yang, and K. Q, "False rumors detection on sina weibo by propagation structures," IEEE Int. Conf. Data Engineering, 2015.

**Tallapudi Jyothsna/** International Journal of Management Research & Review

[7] S. Sun, H. Liu, J. He, and X. Du, "Detecting event rumors on sina weibo automatically," Web Technologies and Applications, Springer, pp. 120–131, 2013.

[8] Zhiwei Jin, Juan Cao,Yongdong Zhang, Jianshe Zhou, and Qi Tian Fellow, "Novel Visual and Statistical Image Features for Microblogs News Verification," IEEE Trans. Inf. Multimedia, pp. 1520-9210, 2016.

[9] Sanjay Yadav and Sanyam Shukla, "Analysis of k-Fold CrossValidation over Hold-Out Validation on Colossal Datasets for Quality Classification," IEEE Int. Conf. Advanced Computing, 2016.

[10] Yuanfang Guo, Xiaochun Cao, Wei Zhang and Rui Wang, "Fake Colorized Image Detection," IEEE Trans. Inf. Information forensics and security, pp. 1556-6013, 2018.