

# TWITTER NEWS STRATIFICATION USING RANDOM FOREST

<sup>1</sup>PRASAD J. KOYANDE, <sup>2</sup>KAVITA P. SHIRSAT

<sup>1,2</sup>Computer Engineering, Vidyalankar Institute of Technology, Mumbai, India  
Email: <sup>1</sup>prasad.koyande@vpt.edu.in, <sup>2</sup>kavita.shirsat@vpt.edu.in

**Abstract:** With the popularity of Social Networks, mostly news providers used to share their news in various social networking sites and web blogs. In India, many news groups share their news on Twitter micro blogging service provider. These data carries valuable information relevant to social research areas. Thus, the idea is to categorize the news into different groups so the news groups in India are identified. News groups are selected on their popularity to extract the short messages from Twitter Micro Blog. Short message extracted from Twitter was classified into 12 major groups. Machine learning techniques were used to train the data. In order to create the instances words from each short message were consider and bag-of-words approach was used to create feature vector. The data was trained using Random Forest machine learning techniques. Random forest is a best ensemble learning method, which is consist of multiple decision trees built on random inputs and separating nodes on a random subset of features. Because of its good classification and generalization ability, random forest is preferred in various domains. Large amount of feature will be collected for current research. The performance will speak the efficacious of the system.

**Keywords:** Random Forest, Web Mining, Text classification.

## I. INTRODUCTION

With the development of web blogs and Social Networks, many organizations and News providers used to share their news headlines in various websites and web blogs. Now-a-days in India, many news groups share their news headlines as short messages in microblogging services such as Twitter. Once these short messages get processed, it carries out significant amount of information which is relevant for many social research areas. The idea is to classify news short messages into different groups so that the user identifies the most popular news group in India. This will help in availing information regarding development, war, education etc. thereby understanding the current state of the country.

There are several news portals currently available to retrieve short messages. Twitter microblog was opted due to 4 reasons [1].

- 1) Microblogging platforms are used by diverse personalities to express their opinion about contrasting topics, thus it is a valuable source of people's speculation.
- 2) Twitter contains an enormous number of text posts and it grows every day. The collected corpus can be arbitrarily large.
- 3) Twitter's audience varies from regular users to celebrities, company representatives, politicians and even the country's president. Therefore, it is possible to collect text posts of users from different social and significant groups.
- 4) Twitter's audience is represented by users from many countries.

The research was conducted and several active Twitter news groups such as 'abnewstv', 'IBN7', 'ibnlive', 'ZoomTV', 'ndtv', 'zeenews', 'timesofindia', 'timesnow', 'economictimes', 'starsportsindia', 'dna',

and 'httweets' were chosen to extract the data. Twitter API provides the ability of extracting 180 short messages with other necessary details in XML format. The short messages will be classified by the system into 12 groups: war-terrorist-crime, economy-business, health, sports, development-government, politics, accident, entertain, disaster-climate, education, society and international. These groups were designated in order to cover the main areas of general news provider.

With the development of machine learning techniques [2], now-a-days, many researchers tend to use machine learning techniques in text classification [1]. There are 2 types of machine learning techniques as supervised learning [3] (the learning data will be provided by the developer) and unsupervised learning [3] (the method will learn a clustering procedure by observing the distance among data). For the present study, supervised learning techniques will be used as the 12 groups do not change regularly.

In order to classify short messages using machine learning techniques, a proper set of features are required to extract from the short messages. The bag-of-words approach [1], was used to extract features from the short messages. The frequency of each word had been used as data. As there are large amount of words in different short messages, using all data will cause to increase the dimension. Thus, the common words are needed to identify and remove from the dataset. Zipf's law [5] states that the utmost common word in a human language text occurs with a frequency inversely proportional to  $n$ . Rarely occurred words do not carry sufficient information but the noise. Thus, low frequent words and high frequent words will be removed from the data in order to reduce the dimension up to certain level.

Once created the dataset, it is important to find a suitable classification method in order to classify

the short messages. Random Forest [6] was used to classify the data as it is capable of dealing with high dimensional dataset [1], [5]. The system provides an accuracy of identifying news which belongs to group accident with 100%, development-government with 100%, disaster-climate with 100%, economy-business with 100%, education with 96.4%, entertainment with 100%, health with 95.2%, international with 80%, politics with 100%, society with 78.40%, sports with 100% and war-terrorist-crime with 100%. The next section will brief out the method of data gathering and feature selection. Section 3 will throw a light on the approach of data training and Section 4 will brief out the evaluation criteria of the system. The general discussion will be brief in section 5.

## II. DATA GATHERING AND FEATURE SELECTION

The classification will be applied into the short messages-news of Twitter microblog. Thus, twitter short messages are needed to be collected. For twitter, there was a character limit such that the length of one short message was limited to 140 characters [7]. Thus the user bound to provide the news by using few amount of words. This caused to limit the words of the short messages into key words. Twitter API provides the ability of retrieving such short messages for a given user in XML file format.

Each XML file could carry out 180 short messages at once. Once gathered the data, the features are need to extract from the short messages. These features are required to learn the patterns amount the groups. The words are used as features. Thus the bag-of-words approach [1], [5] was used to extract the features. This will pool the words from all short messages and will create a document vector, containing words. Some researchers had used n-gram instead of words [8]. However, n-gram method cause to increase the dimension of the dataset, as it uses unigram, bigram, trigram which make complex for the system to recognize the pattern [8]. Thus, the words are chosen as the features.

In order to create the dataset, the frequencies of words were used. All words of the documents do not carry out useful information. To avoid using very low frequent words which do not carry out any valuable information regarding the group. Zipf's law [5] states that the  $n$ th most common word in a human language text occurs with a frequency inversely proportional to  $n$ . Thus, the common words were removed from the dataset by removing high frequency data. Therefore, a lower cut off value and upper cut off value were required to choose in order to obtain best set of features. The frequencies of selected data range were shown in figure 1. The values  $X$  and  $Y$  were chosen as the frequency limit (lower cut off value and upper cut off value) which maximize the effectiveness.

Thus, unimportant features were removed from the dataset. This caused to reduce the dimension which makes the training method more effective.



FIGURE 1: RANGE OF FREQUENCIES OF CHOSEN DATA

In machine learning [2], [3], there are basically two types of learning methods. Supervised learning [2] and unsupervised learning [2]. In supervised learning, the developer provides learning data to the system in order to train the system. In unsupervised learning, the system itself learns patterns from the data. For the current situation, as the group are predefined and do not change regularly, supervised learning method is more applicable [3].

## III. DATA TRAINING

In machine learning [2], [3], there are basically two types of learning methods. Supervised learning [2] and unsupervised learning [2]. In supervised learning, the developer provides learning data to the system in order to train the system. In unsupervised learning, the system itself learns patterns from the data. For the current situation, as the group are pre-defined and do not change regularly, supervised learning method is more applicable [3].

Ensemble learning [6] refers to the algorithms that produce collections or ensembles of classifiers which learn to classify by training individual learners and fusing their predictions. Growing an ensemble of trees and getting them vote for the most popular class has provided a good enhancement in the accuracy of classification. Often, random vectors are built that control the growth of each tree in the ensemble. The ensemble learning methods can be divided into two main groups: bagging and boosting. In bagging, models are fit in parallel where successive trees do not depend on previous trees. Each tree is independently built using bootstrap sample of the dataset. A majority vote determines prediction. In boosting, models are fit sequentially where successive trees assign additional weight to those observations poorly predicted by previous model. A weighted vote specifies prediction.

A random forest [6] adds an additional degree of randomness to bagging. Although each tree is constructed using a different bootstrap sample of the dataset, the method by which the classification trees are built is improved. A random forest predictor is an ensemble of individual classification tree predictors. For each observation, each individual tree votes for one class and the forest predicts the class that has the plurality of votes. The user has to specify the number of randomly selected variables ( $m_{try}$ ) to be searched through for the best split at each node.

Whilst a node is split using the best split among all variables in standard trees, in a random forest the node is split using the best among a subset of predictors randomly chosen at that node. The largest tree possible is grown and is not pruned. The root node of each tree in the forest contains a bootstrap sample from the original data as the training set. The observations that are not in the training set, are referred to as "out-of-bag" observations.

Since an individual tree is unpruned, the terminal nodes can contain only a small number of observations. The training data are run down each tree. If observations  $i$  and  $j$  both end up in the same terminal node, the similarity between  $i$  and  $j$  is increased by one. At the end of the forest construction, the similarities are symmetrized and divided by the number of trees. The similarity between an observation and itself is set to one. The similarities between objects form a matrix which is symmetric, and each entry lies in the unit interval  $[0,1]$ . Breiman defines the random forest as [6]:

A random forest is a classifier consisting of a collection of tree-structured classifiers  $\{h(x, e_k), k = 1, \dots\}$  where  $\{e_k\}$  are independent identically distributed random vectors and each tree casts a unit vote for the most popular class at input  $x$ .

A summary of the random forest algorithm for classification is given below [9]:

- Draw  $n_{tree}$  bootstrap samples from the original data.
- For each of the bootstrap samples, grow an unpruned classification tree, with the following modification: at each node, rather than choosing the best split among all predictors, randomly sample  $m_{try}$  of the predictors and choose the best split from among those variables. Bagging can be thought of as the special case of the random forest obtained when  $m_{try} = p$ , the number of predictors.
- Predict new data by aggregating the predictions of the  $n_{tree}$  trees, i.e., majority votes for classification, average for regression.

The generalization error of a forest of tree classifiers depends on the strength of the individual trees in the forest and the correlation between them. Using a random selection of features to split each node yields error rates that compare to AdaBoost [9]. An estimate of the error rate can be obtained, based on the training data, by the following [9]:

- At each bootstrap iteration, predict the data that is not in the bootstrap sample, called "out-of-bag" data, using the tree which is grown with the bootstrap sample.
- Aggregate the out-of-bag predictions. On the average, each data point would be out-of-bag around 36% of the times, so aggregate these

predictions. Calculate the error rate, and call it the "out-of-bag" estimate of error rate.

The random forest performs well compared to several other popular classifiers. In addition, it is user friendly as it has only two parameters: (i) the number of variables in the random subset at each node, and (ii) the number of trees in the forest. The random forest is not usually very sensitive to the values of these parameters.

To create the training data and testing data, each short message was classified to a group manually. One short message might belong into several groups. Therefore, each category was considered as a separate binary classification problem [4]. The training process was developed in order to recognize whether the selected short message belongs to the group A short messages will be classified manually as "Group A" or "other". 90% data was used to train the system and 10% were used to test the system [4]. In order to train the 12 groups, there will be 12 training data sets and testing will result in 12 tables, each table describing their performance as in table 1.

TABLE 1

Observed classes	Expected Group	
	True positive (tp)	False positive (fp)
False negative (fn)	True negative (tn)	

#### IV. EVALUATION

The evaluation was carried out in order to measure the effectiveness. Effectiveness is purely a measure of the ability of the system to satisfy the user in terms of the relevance of short messages retrieved [12]. It is assumed that the more effective the system, the more it will satisfy the user [12]. The effectiveness of the retrieval system was measured using precision and recall values [12]. Precision is the fraction of retrieved short messages that are relevant. Recall is the fraction of relevant short messages that are retrieved [13]. As the system results the performance as in Table 1, the precision can be calculated using Equation 1 and recall could be calculated using Equation 2.

$$Precision = \frac{t_p}{t_p + f_p} \quad (1)$$

$$Recall = \frac{t_p}{t_p + f_n} \quad (2)$$

In order to measure the performance, the performance will be depending on the biasness of the training data. In order to avoid the biasness, cross validation [14] was applied to the testing process. Koveri [14] suggest that the best number of fold is 10. Thus, 90% data were randomly chosen to train the data. The system was tested 10 times per group and the average precision and

recall values were recalculated [12]. The results were briefed in table 2.

These measures were used to figure out best frequency limit for the feature selection. Thus, it is important to calculate a single measurement instead of 2 values [14]. Many alternative methods were proposed over the years and Harmonic mean [13] had identified as the best single value summaries [13]. The harmonic mean (F-measure) was given in Equation 3 and weighted harmonic mean was given in Equation 4. Values of P less than 1 emphasize the precision whereas values of P greater than 1 emphasizes recall [13].

$$F = \frac{2 \times P \times R}{P + R} \quad (3)$$

$$F_{\beta} = \frac{(\beta^2 + 1) \times P \times R}{\beta^2 \times P + R} \quad (4)$$

For the current situation, the main focus will be on the fraction of retrieved short messages that are relevant. Thus  $F_{0.5}$  will be used to select the best frequency range and to measure the effectiveness of the system. Table 3 provides the values obtained for F measure and  $F_{0.5}$ .

TABLE 2: PRECISION AND RECALL VALUES

Group Name	Avg. Precision	Avg. Recall
Accident	0.824	1
Development-Government	1	1
Climate-Disaster	0.733	1
Economy-Business	0.955	1
Education	0.964	0.964
Entertainment	0.846	1
Health	0.714	0.952
International	1	0.8
Politics	0.8	1
Society	0.965	0.784
Sports	0.857	1
War-Terrorist-Crime	1	1

TABLE 3: RESULTS OF F-MEASURE AND  $F_{\beta}$ -MEASURE

Group Name	F-MEASURE	$F_{\beta}$ -MEASURE
Accident	0.903	1
Development-Government	1	1
Climate-Disaster	0.846	1
Economy-Business	0.977	1
Education	0.964	1
Entertainment	0.917	1
Health	0.816	0.994
International	0.889	0.999
Politics	0.889	1
Society	0.865	0.961
Sports	0.923	1
War-Terrorist-Crime	1	1

With the results obtained from Table 3 it is clear that Random Forest provides good results for most of groups.

## V. DISCUSSION

A system which is able to classify news headlines will be useful in various social researches. With the

development of web technologies, people get involved in many social networks and web blogs. Twitter is a micro blog which allows many famous news suppliers to publish their news headlines. Twitter API supports user to retrieve available short messages. These retrieved files will be in XML file format and each file could retrieve maximum number of 200 short messages per once.

In order to apply machine learning, a proper feature set was required. The feature set was created by pooling the words and creating a document vector. This approach was named as bag-of-words approach. The frequency of each word was chosen as data. When considering all words together, it creates a huge dimension of instances. In order to reduce the dimension, a lower cut off frequency and an upper cut off frequency value were chosen. The value was chosen as the frequency range which maximizes the accuracy.

There are 12 groups defined and each group was treated as separate binary classification problem as same short message could belong into several groups. System was trained using Random Forest. The effectiveness of the training system can be measured using recall and precision values. Precision is the probability of retrieving relevant short messages. Recall is the probability of the relevancy of retrieved short messages. The harmonic measure (F-measure) was used to obtain a single value for recall and precision. The weighted F-measure ( $F_{\beta}$  measure) was used as precision was needed to be emphasized in current situation. The system provides best results for Accident, Development-Government, Climate-Disaster, Entertainment, Health, Education, Sports, War-Terrorist-Crime, Politics and Economy-business groups. It provides reasonable effectiveness more than 96%.

## REFERENCES

- [1] A. Pak and P. Paroubek, "Twitter as a Corpus for Sentiment Analysis and Opinion," in *Analysis*, 2010.
- [2] N. J. Nilsson, *INTRODUCTION TO MACHINE LEARNING*. 1998.
- [3] J. K. M. Han, *Data Mining : Concepts and Techniques*, 2<sup>nd</sup> ed. 2006.
- [4] Inoshika Dilrukshi, Kasun De Zoysa, Amitha Caldera. "Twitter News Classification Using SVM", *Computer Science & Education (ICCSE 2013 IEEE)*.
- [5] K. G. Zipf, *Human Behaviour and the Principle of Least Effort*. Oxford, England: Addison-Wesley, 1949.
- [6] L. Breiman, "Random Forests," *Machine Learning*, Vol. 40, No. 1, 2001.
- [7] (2012, Apr.) Counting Characters. [Online]. <https://dev.twitter.com/docs/counting-characters>
- [8] W. B. Cavnar and T. J. M., "N-Gram-Based Text Categorization," in *Proceedings of SDAIR-94, 3rd Annual Symposium on Document Analysis and Information Retrieval*, pp. 161-175, 1994.
- [9] Andy Liaw and Matthew Wiener, "Classification and Regression by Random Forest," *R News*, Vol. 2/3, December 2002.

- [10] Yoav Freund and Robert E. Schapire, "A Short Introduction to Boosting", Journal of Japanese Society for Artificial Intelligence, 14(5):771-780, September, 1999.
- [11] A.Z. Kouzani, S. Nahavandi, K. Khoshmanesh,"Face Classification by a Random Forest",TENCON 2007 - 2007 IEEE Region 10 Conference,Oct. 30 2007-Nov. 2 2007.
- [12] C. I Rijsbergen, Information Retrieval, 2nd ed. London: Butterworths, 1979.
- [13] C. D. Manning, P. Raghavan, and H. Schutze, Introduction to Information Retrieval.Cambridge University Press, 2008.
- [14] Y. Baeza and B. R. Neto, Modern Information Retrieval.Boston, 1999.

★ ★ ★