

Enhancing Spam Detection on Nepali Language SMS

Anjana Prajapati

Nepal College of Information Technology
Pokhara University, Nepal
anjanapasa2004@gmail.com

Arun Timalsina*

Pulchowk Campus, IoE,
Tribhuvan University, Nepal
t.arun@ieee.org

*Corresponding author

Abstract - Any junk message that is delivered to a mobile phone as text messaging through the SMS is called SMS spam. As the popularity of mobile phone devices has increased over the recent years, SMS has grown into a multi-billion dollar industry. At the same time, the reduction in the cost of messaging services has resulted in growth in unsolicited commercial advertisements (spam) being sent to mobile phones. The most common filtering technique is content-based filtering which uses the actual text of the message to determine whether it is spam or ham. Since the characteristics used by the filter to identify spam message are constantly changing over time, it is very challenging to represent all information in a mathematical model of classification. The Nepali language is morphologically rich and it is a challenging task to build a model for such language. Different supervised learning classifiers, Decision Tree, Logistic Regression, NB and SVM along with the combination of SVM and NB classifiers as SVM-NB are used for classification of spam and ham on Nepali text in this research. The accuracy, recall rate, precision rate and f1-score for these classifiers are analyzed.

Keywords: spam filtering, spam, ham, Decision Tree, Logistic Regression, Naive Bayes, Support Vector Machine.

I. INTRODUCTION

The mobile phone market has experienced substantial growth over recent years. It is the most widely used data application and its popularity has increased, there is a number of unsolicited commercial advertisements sent to mobile phones using text messaging. SMS text messages are indispensable in our lives today, but along with the convenience of using SMS messages in our daily lives, we also face a serious problem caused by SMS spamming.

Spamming is the abuse of electronic messaging systems to send unsolicited bulk messages or to promote products or services, which are almost universally, undesired bulk data received by massive recipients.

Spamming is a serious problem for SMS today, as it is for emails and social networking services, because it disrupts people's daily life and harms the well-being of telecom operators. SMS has certain characters that are different from emails. So there exists a difference between spam-filtering in text messages and emails. Emails have a variety of large datasets available whereas real databases for SMS spams are very limited. Additionally, the number of features that can be used for their classification is far smaller than the corresponding number in emails due to the small length of text messages.

The mobile technology has been even growing in Nepal too. People nowadays want to use Nepali Unicode for communication in SMS. Now and then, we also receive spam SMS in Nepali language promoting their business. There are a lot of spam filters for the English language already available, but for Nepali text, it is very rare. So, the spam detection system for the Nepali language can help telecom operators to filter out such SMS from even being sent to the end users.

Several researchers have applied machine learning techniques in order to improve the detection of spam messages. Spam filtering is the processing of data to organize it according to specified criteria. There are many spam filters available and the motivation behind developing spam filters is due to its wide range of applications like mobile application and computer system.

II. RELATED WORKS

There has been significant amount of research done in the field of spam filtering in case of English Language. But we can find very few such researches on Nepali text. We can refer to similar techniques used for English spam classification in Nepali text spam classification. However, preprocessing involves significantly different steps for Nepali language. Various researchers have provided different way of combining the different machine learning classifiers for spam

Diksha S. Jawale, Ashwini G. Mahajan, Kalyani R. Shinkar and Vaishnavi V. Katdare. [1] and Weimiao Feng, Qing Yang, Jianguo Sung, Liguozhang, Cuiling Cao [2] have combined the SVM and NB method for email classification technique and claimed the accuracy as well as execution time improvement with this combination. The former named the method as NB-SVM. The later named the combination as SVM-NB. In [1], the research was carried out on dataset from datamall.

In case of Nepal text spam classification, very few research work can be found online. Tej Bahadur Shahi and Abhimanu Yadav [3] have used NB and SVM individually for the spam classification for mobile SMS with Nepali Text. The calculated accuracy for SVM was 87.15% and 92.74% with their dataset. Similarly, Heramba Raj Pant [4] has used SVM and Multilayer Neural network for Nepali text SMS classification.

NB classifiers and SVM classifiers are most popular classifiers used for email classification. For both SMS and email, since the concern is content based classification, the algorithms would work for both in similar way. There are number of forms of Naïve Bayes that have been used for email classifications. Significant amount of work has been done using Naïve Bayes in different forms [5], [6], [7], [8] etc. Similarly, there has been significant research on SVM classifiers [9], [10].

Saumya Goyal, Prof. R. K. Chauhan, Shabnam Parveen [11] have used KNN and Decision Tree mechanism in

social network to detect spam and normal messages and links in their research paper. Weka tool was used to analyze the algorithms on real dataset and they resulted that KNN was better than Decision Tree.

Alhassan J. Ibrahim, Maheyzah Md Siraj, Mazura Mat Din [12] have also carried out research on spam review detection. They have used NB, SVM and LR in their research. They have used WEKA to implement the ensemble classification phase which was possible through the choice voting algorithm provided and claimed that the ensemble classification provided better results.

III. RESEARCH FRAMEWORK

The purpose of this research is to enhance the spam detection in context of Nepali SMS and to analyze the performance of different supervised machine learning models for the Nepali SMS spam detection. Figure 1 shows the summary of the research framework used in terms of block diagram. The SMS collected from difference sources are selected for the experimentation. The collected data is preprocessed with stemming and stopwords removal. The data is tokenized as bag of words and the feature selection is done. The total collected data is split into training set and test set. K-Fold validation has been used to split the data into training and test set. The train data is trained with various machine learning models. The test data is used for the prediction and the results are obtained. The results obtained from various models are evaluated and compared with each other.

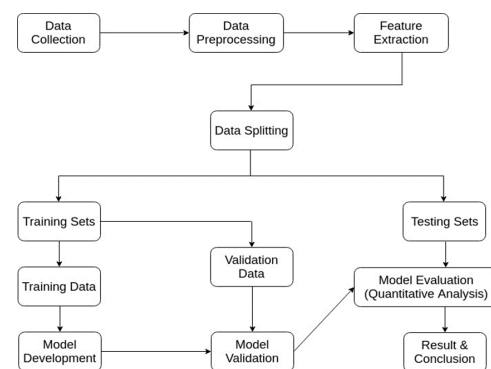


Figure 1 Research Framework

A. Preprocessing

Figure 2 shows the preprocessing steps applied for the collected dataset of 500 SMS with 339 ham and 161 spam messages arranged in random order.

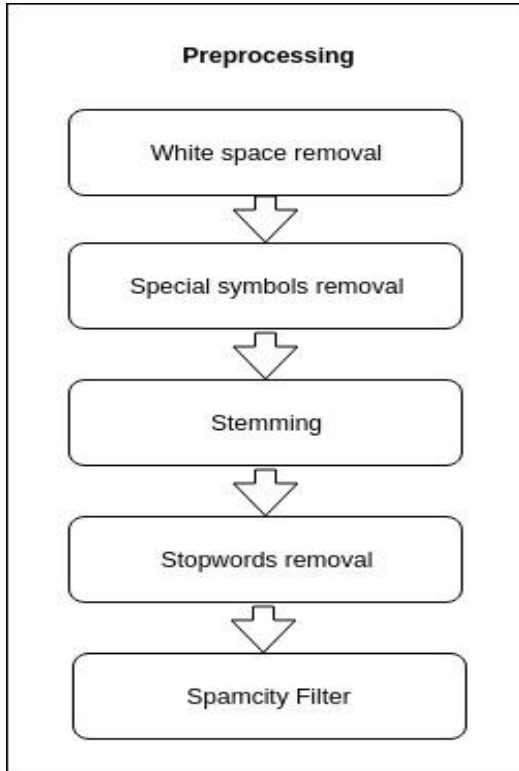


Figure 2 Preprocessing

B. Spamicity

The concept of spamicity has been used to filter such the words which can contribute less to the classification. Spamicity determines the likeliness of a particular word to be spam. This likeliness is calculated in terms of probability.

If,

$P(w/S)$ = probability that the specific word appears in spam message,

$P(w/H)$ = probability that the specific word appears in ham message,

The spamicity of the word is mathematically defined as,

$$S(w) = P(w/S) \cdot P(w/S) + P(w/H) \tag{1}$$

If spamicity is close to 1, it is a good spam indicator and if spamicity is close to 0, it is good ham indicator. Hence high rank is assigned to words which are far away from 0.5.

However, there are also the cases where both $P(w/S)$ and $P(w/H)$ is close to too small and close to 0, which cannot be good spam indicator even if it is far from 0.5. Hence the absolute difference $|P(w/S) - P(w/H)|$ is calculated to sort the words.

Hence the algorithm for feature selection would be:

1. Filter out words with $|spamicity - 0.5| < 0.05$
2. For each remaining word calculate $|P(w/S) - P(w/H)|$
3. Sort the words as per value calculated in step 2 and select top k words.

The preprocessing involves filtering out the words based upon spamicity as shown in step 1.

IV. EXPERIMENT AND ANALYSIS

The data sample consists of the 500 SMS messages with 339 spam messages and 161 spam messages which are all arranged in random order. The messages are exactly the same way they are received in the mobile devices in raw format without preprocessing. The data frame is created from these list of messages along with its spam and ham labels saved in a CSV file.

The accuracy, precision, recall and f1-score are calculated for all five methods with 10-fold validation technique.

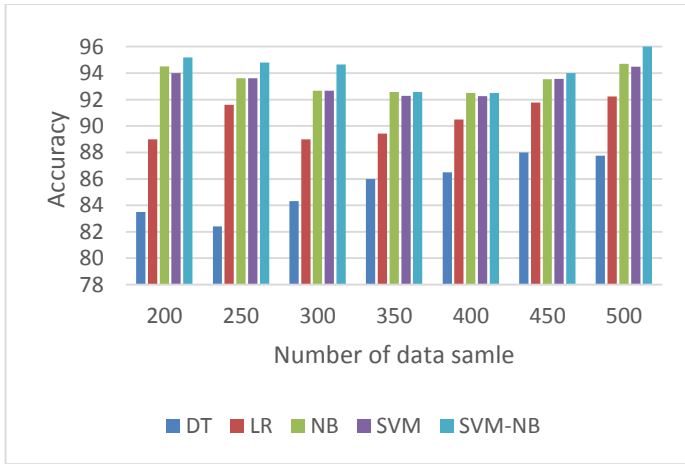


Figure 3 Accuracy for k=10

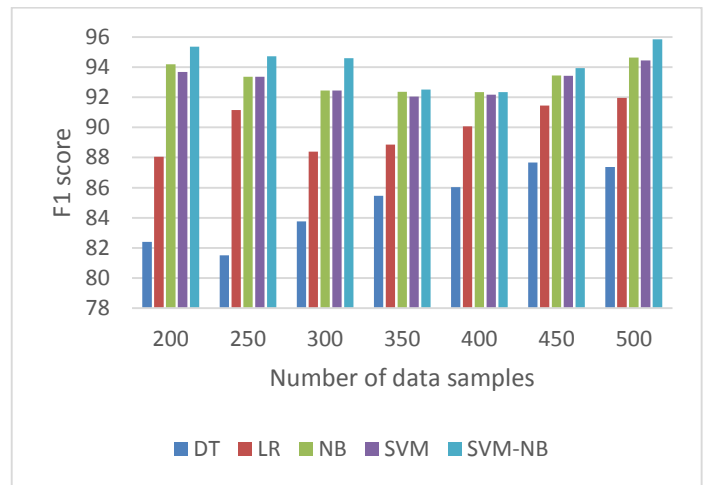


Figure 6 f1-score for k=10

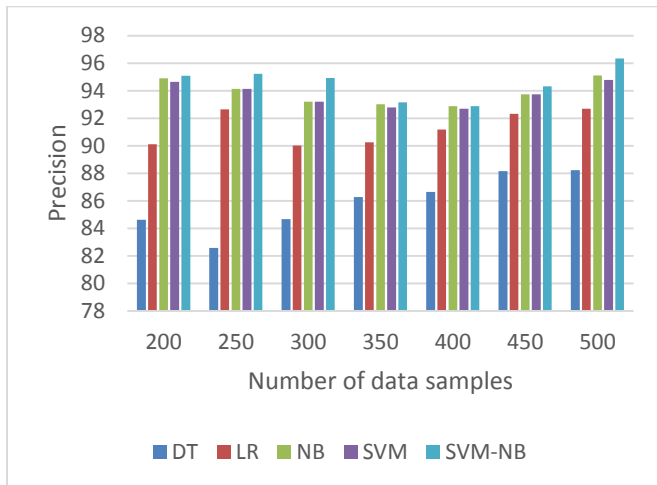


Figure 4 Precision for k=10

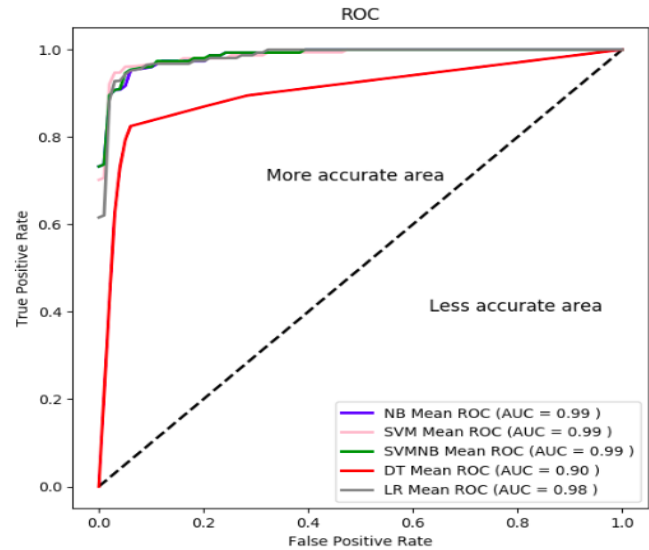


Figure 7 Mean ROC curve

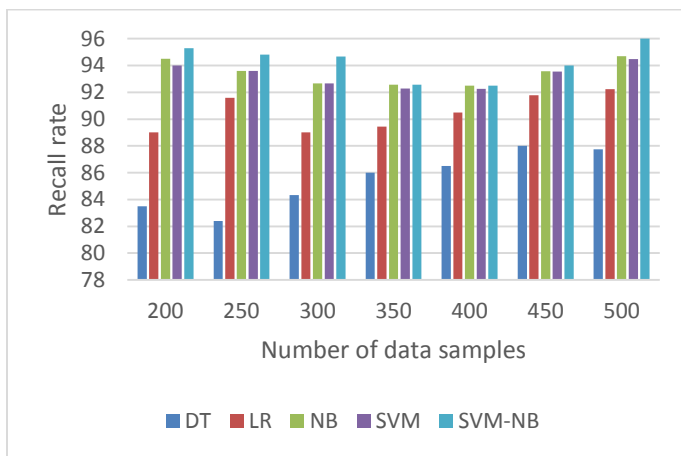


Figure 5 Recall rate for k=10

The experiment and result analysis figures show the result of performance measures for the all five models calculated by providing the number of data samples. From the above result, it is observed that SVM-NB performs better result among all other methods.

V. CONCLUSION

A combined model of SVM and NB as SVM-NB along with other four individual methods Decision Tree, Logistic Regression, NB and SVM has been used for the spam detection in Nepali SMS in this thesis. Among all the five models used, SVM-NB resulted with better

performance. The performance was compared based on the metrics accuracy, precision, recall and f1-score.

In paper [3] individual NB and SVM were used for Nepali SMS spam detection with accuracy 92.74% in NB and 87.15% in SVM. The hybrid method SVM-NB was implemented in the papers [1] and [2] but for English email dataset and Chinese email dataset respectively. In this thesis, the hybrid SVM-NB method was implemented using spamicity for Nepali language based SMS spam detection.

With 10-fold validation technique, SVM-NB was found to have highest accuracy and precision with 96.00% accuracy and 96.35% precision whereas, Decision Tree was found to be least effective with the categorization with 87.75% accuracy and 88.22% precision for 500 data samples.

VI. FUTURE WORK

The future works to improve the Nepali SMS spam detection can be done in a number of ways. To mention a few of them, the comparison with other machine learning models like Neural Network can be done. The collection of data was limited to 500 SMS data including 339 ham SMS and 161 spam SMS. It can be improved by taking more data samples into consideration. The research uses rule-based stemming in the basic form. More effective use of stemming can help to improve the overall accuracy of the classification. The numbers used to send the SMS are not considered in the dataset. Considering the phone number as a feature may improve the accuracy of the spam detection.

ACKNOWLEDGMENT

This research is supported by the Nepal College of Information Technology (NCIT), Department of Computer Science and Engineering, Pokhara University, Nepal.

REFERENCES

- [1] D. S. Jawale, A. G. Mahajan, K. R. Shinkar and V. V. Katdare, "Hybrid spam detection using machine learning," *International Journal of Advance Research, Ideas and Innovations in Technology*, vol. 4, no. 2, pp. 2828-2832, 2018.
- [2] W. Feng, J. Sun, L. Zhang, C. Cao and Q. Yang, "A Support Vector Machine based Naive Bayes Algorithm for Spam Filtering," in *2016 IEEE 35th International Performance Computing and Communications Conference (IPCCC)*, Las Vegas, NV, USA, 2016.
- [3] T. B. Shahi and A. Yadav, "Mobile SMS Spam Filtering for Nepali Text Using Naive Bayesian and Support Vector Machine," *International Journal of Intelligence Science*, vol. 4, pp. 24-28, 2014.
- [4] H. R. Pant, "Mobile SMS Spam Filtering For Nepali Text Using Neural Network and Support Vector Machine Classifier".
- [5] H. Zhang and D. Li, "Naïve Bayes Text Classifier," in *IEEE International Conference on Granular Computing (GRC 2007)*, Fremont, CA, USA, 2007.
- [6] Z. Yang, X. Nie, W. Xu and J. Guo, "An Approach to Spam Detection by Naive Bayes Ensemble Based on Decision Induction," in *Sixth International Conference on Intelligent Systems Design and Applications*, Jinan, China, 2006.
- [7] S.-B. Kim, K.-S. Han, H.-C. Rim and S. H. Myaeng, "Some Effective Techniques for Naive Bayes Text Classification," *IEEE Transactions on Knowledge and Data Engineering*, vol. 18, no. 11, pp. 1457-1466, 2006.
- [8] A. K. Seewald, "An evaluation of Naive Bayes variants in content-based learning for spam filtering," *Intelligent Data Analysis*, vol. 11, no. 5, pp. 497-524, 2007.
- [9] H. Drucker and D. Wu, "Support Vector Machines for Spam Categorization," *IEEE Transactions on Neural Networks*, vol. 10, no. 5, pp. 1048-1054, 1999.
- [10] S. and T. Verma, "E-Mail Spam Detection and Classification Using SVM and Feature Extraction," *International Journal of Advance Research, Ideas and Innovations in Technology*, vol. 3, no. 3, pp. 1491-1495, 2017.
- [11] S. Goyal, R. K. Chauhan and S. Parveen, "Spam detect using KNN and decision tree mechanism in social network," in *Fourth International Conference on Parallel, Distributed and Grid Computing (PDGC)*, Wanknaghat, India, 2016.
- [12] A. J. Ibrahim, M. M. Siraj and M. M. Din, "Ensemble classifiers for spam review detection," in *IEEE Conference on Application, Information and Network Security (AINS)*, Miri, Malaysia, 2017.
- [13] S. K. Trivedi and S. Dey, "A combining Classifiers Approach for detecting Email Spams," in *Advanced Information Networking and Application Workshops (WAINA)*, Crans-Montana, 2016.
- [14] N. F. Rusland, N. Wahid, S. Kasim and H. Hafit, "Analysis of Naïve Bayes Algorithm for Email Spam Filtering across Multiple Datasets," in *International Research and Innovation Summit (IRIS2017)*, Melaka,

Malaysia, 2017.

- [15] H. Sajedi, G. Z. Parast and F. Akbari, "SMS Spam Filtering Using Machine Learning Techniques: A Survey," *Machine Learning Research*, vol. 1, no. 1, pp. 1-14, 2016.
- [16] P. Navaney, G. Dubey and A. Rana, "SMS Spam Filtering Using Supervised Machine Learning Algorithms," in *8th International Conference on Cloud Computing, Data Science & Engineering (Confluence)*, Noida, India, 2018.