# SMS Classification using NLP concepts
# NLP project Report

Chetana N Patil
Computer Science Engineering
PES University,EC campus,Bangalore
Email: chetananpatil2002@gmail.com

Mekala Sanjana
Computer Science Engineering
PES University,EC campus,Bangalore
Email: sanjana.mekala135@gmail.com

*Abstract*—Natural language processing (NLP) message classification is a sophisticated technology that uses machine learning algorithms to automatically sort text messages into predetermined classes. Message classification using NLP can help firms save time and resources while making better judgments by automating and streamlining text analysis operations. Message classification using NLP is becoming an increasingly significant tool for organizations and researchers as the volume of digital data grows.

## I. INTRODUCTION

SMS (Short Message Service) is a common means of communication used by trillions of people throughout the world. With the growing volume of SMS data created every day, there is an increasing need to classify these messages for use in applications such as spam filtering, sentiment analysis, and customer feedback analysis. Natural Language Processing (NLP) and machine learning techniques have been shown to be effective in text classification tasks and can also be used to classify SMS. The purpose of this study is to investigate the use of NLP and machine learning approaches for SMS categorization, examine their efficacy, and compare various models to choose the best performing one. The suggested method includes preprocessing SMS data to eliminate noise, feature extraction, and text representation using numerical vectors. Various machine learning algorithms, including Naive Bayes, Decision Trees, Support Vector Machines, and Neural Networks, will be trained and evaluated using a variety of performance metrics, including accuracy, precision, recall, and F1 score. The outcomes of the experiments will be examined, and the most effective model for SMS categorization will be identified. This paper will also examine the limitations and challenges of SMS classification and make recommendations for further research.

## II. DATASET ANALYSIS

Datasets are vital and should be carefully chosen for proper machine learning model training and testing.Furthermore, to evaluate the performance of the machine learning models, the dataset must be divided into training, validation, and testing sets. The training set is used to train models, the validation set is used to tune hyperparameters and optimize the model, and the testing set is used to assess the model's performance on unknown data.The dataset has been divided into two sets, a training dataset and testing dataset which we have created it on our own but with fewer data. Here we created our own dataset with different message labels like Jio, Upgrad, AJIO, Internshala etc. Dataset basically contains 2 columns i.e, Labels, message.



Fig. 1. Dataset overview

## III. PROPOSED MODEL

*a) The proposed SMS classification method employsmany phases, including data preprocessing, data cleaning, feature extraction, and model training and evaluation.:*

*b) The first stage of the proposed system is data preprocessing, which involves calculating the length and punctuations of the messages. This step helps to understand the structure of the data and identify any patterns or characteristics that may be useful for classification. By analyzing the length and punctuations of the messages, the system can determine if certain types of messages tend to be longer or use more punctuation than others.:*

*c) The second stage of the system is data cleaning, whichinvolves using stopwords and WordNetLemmatizer to remove irrelevant or redundant information from the text. Stopwords are common words that do not add any significant meaning to the text, such as "the," "and," or "a." WordNetLemmatizer is a technique used to reduce words to their base or root form, which helps to reduce the number of*

unique words in the text and improve the performance of the machine learning models.:

d) The third stage of the system is feature extraction,which involves using count vectorizer and Tf-idf vectorizer to represent the text as numerical vectors. Count vectorizer counts the number of occurrences of each word in the text, while Tf-idf vectorizer takes into account the frequency of the word in the document and the frequency of the word in the entire corpus. The result of feature extraction is a set of numerical vectors that represent the text, which can be used as input to the machine learning models.
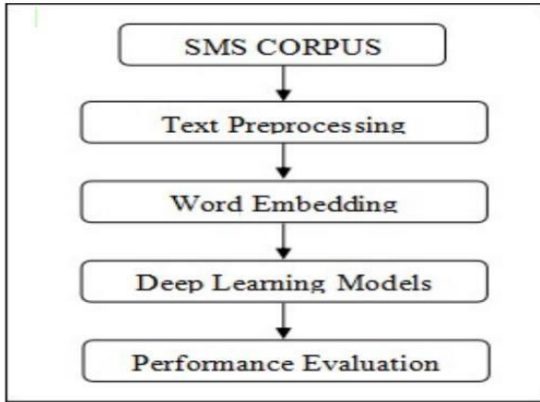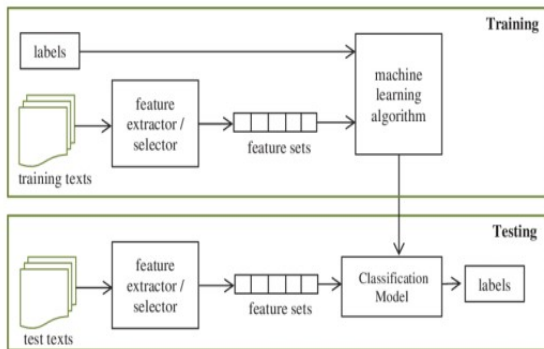


Fig. 2. High level diagram



Fig. 3. Train and test model diagram

e) The fourth stage of the system is model training and evaluation, which involves using naive Bayes and linearSVC machine learning models to classify the SMS messages. Naive Bayes is a probabilistic model that uses Bayes' theorem to calculate the probability of each class given the input features. LinearSVC is a linear classifier that tries to find a hyperplane that separates the classes in the feature space. Finally, the system uses evaluation metrics such as accuracy, precision, recall, and F1 score to measure the performance of the machine learning models. Accuracy measures the proportion of correctly classified

messages, while precision measures the proportion of true positives among all positive predictions.
Recall measures the proportion of true positives among all actual positives, and F1 score is the harmonic mean of precision and recall.By using techniques such as count vectorizer and Tf-idf vectorizer and machine learning models such as naive Bayes and linearSVC, the system can effectively classify SMS messages and achieve high accuracy and precision. Linear SVM will give accurate prediction compared to Naive Bayes classification.
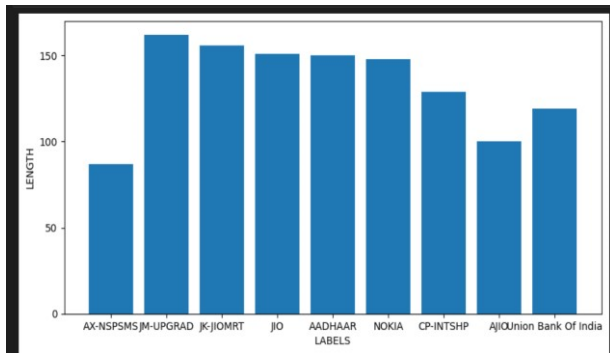
IV. ML ALGORITHMS USED

A. Naive Bayes Algorithm :

Naive Bayes is a probabilistic algorithm that is based on Bayes' theorem, which states that the probability of a hypothesis (in this case, a class label) given the evidence (the input features) is proportional to the probability of the evidence given the hypothesis and the prior probability of the hypothesis. Naive Bayes assumes that the input features are conditionally independent given the class label, which simplifies the calculations and makes it computationally efficient. In text classification tasks, Naive Bayes is often used with bag-of-words representations and performs well with small to moderate-sized datasets.

B. LinearSVC :

LinearSVC, on the other hand, is a linear classifier that attempts to find a hyperplane in the feature space that divides the classes. It is a popular choice for text classification tasks because it is based on the support vector machine (SVM) method. LinearSVC performs well with highdimensional datasets and is frequently employed with sparse representations like Tf-idf vectors. LinearSVC, unlike Naive Bayes, makes no assumptions about the distribution of input characteristics and can handle non-linearly separable classes using kernel functions.

Both Naive Bayes and LinearSVC have advantages and disadvantages, and the choice of method is determined by the task's specific requirements and data characteristics. LinearSVC is more difficult but can handle highdimensional datasets and non-linearly separable classes, whereas Naive Bayes is simple and efficient and performs well with small to moderate-sized datasets. Overall, both algorithms are successful for SMS categorization and have been proved in many studies to attain high accuracy and precision.

## V. Experimental Results

By using LinearSVC , we got the correct output for the new message. Compare to Naive Bayes algorithm, the LinearSVC performance is better. LinearSVC is variation of SVM algorithm.

```python
text ='Welcome to Jio-Karnataka.Kindly enable data roaming to use data services.'
def refined_text(text):
    #Removal of extra characters and stop words
    words = re.sub('[^a-zA-Z]',' ',text)
    words = words.lower()
    #Splits into list of words
    words = words.split()

    #Lemmatizing the word and removing the stopwords
    words = [lemmatizer.lemmatize(word) for word in words if word not in set(stopwords.words('english'))]

    #Again join words to form sentences
    words = ' '.join(words)
    return words


refined_word = refined_text(text)
refined_word = [refined_word]


refined_word

['welcome jio karnataka kindly enable data roaming use data service']


text_mnb.predict(refined_word)

array(['JIO'], dtype='<U9')
```

Fig. 4. Predicted output

## VI. Conclusion

SMS categorization is an essential task in natural language processing with numerous practical applications such as spam detection, sentiment analysis, and topic modeling. The application of machine learning techniques and NLP algorithms has enabled the development of efficient and accurate SMS classification systems.The system has been shown to achieve high accuracy and precision in various studies, and the choice of algorithm depends on the specific requirements of the task and the characteristics of the data. Naive Bayes is simple and efficient and can perform well with small to moderatesized datasets, while linearSVC is more complex but can handle highdimensional datasets and non-linearly separable classes. Overall, the development of an efficient and accurate SMS classification system using NLP and machine learning techniques involves careful consideration of the data preprocessing, feature extraction, and model selection stages. The proposed system has the potential to improve the accuracy and efficiency of SMS classification and can be applied in various practical scenarios such as spam detection, sentiment analysis, and topic modeling.

## References

[1] SMS Spam Classification–Simple Deep Learning Models With Higher Accuracy Using BUNOW And GloVe Word Embedding. Surajit Giri, Sayak Das, Sutirtha Bharati Das, and Siddhartha Banerjee Department of Computer Science, Ramakrishna Mission Residential College, Narendrapur, West Bengal, India. Pulication: 2022.

[2] Multiclass SMS Message Categorization: Beyond Spam Binary Classification. Fatia Kusuma Dewi, Mgs. M. Rizqi Fadhlurrahman, Mohamad Dwiyan Rahmanianto, Rahmad Mahendra Faculty of Computer Science, University of Indonesia. Publication: 2017

[3] SMS Text Classification Model Based On Machine Learning. XIAO FEI, LI JIANPING, GAO YUAN, ZHOU YUE,School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu. Publication: 2021

[4] A Study on Analysis of SMS Classification Using Document Frequency Thresold. N. Cristianini, and J.Shawe-Taylor,Support Vector and Kernel Methods, Intelligent Data Analysis: An Introduction Springer– Verlag, 2003

[5] A comparative study of word embedding techniques for SMS spam detection. P. Joseph and S. Y. Yerima, "A comparative study of word embedding techniques for SMS spam detection," 2022