

Context-based Sentiment analysis of Indian Marathi Text using Deep Learning

¹Ms. Kirti Kakde, ²Dr. H. M. Padalikar

¹Research Scholer: Dept. of Computer Application,
IMED, Bharati Vidyapeeth Deemed to be University,
Pune, India

kirtikakde2007@gmail.com

²Professor, Dept. of Computer Application,
IMED, Bharati Vidyapeeth Deemed to be University,
Pune, India

hmpadalikar1965@gmail.com

Abstract— In Digital India, the Internet plays a crucial role in communication. The English language is widely used for such a process. The Internet has no language barrier. India is a multi-lingual country with boundless linguistic and social diversities. The most trending pattern observed in India is people intend to post their views, thoughts, feedback, and comments in their mother tongue over social media and blogs. Views posted by people is important for organization belonging to any category small, medium and large enterprises to improve their product or service. This data is hastily accumulated every day which should be necessary to identify and process. In terms of processing little work has been done for Indian languages where traditional approaches were used which are far away from the context of the text. In this research to perform sentiment analysis supervised algorithms that is Multinomial Naïve Bayes is implemented on the Marathi dataset. Along with this deep learning, Natural Language Processing approach Bidirectional Encoder Representations from Transformers (BERT) is utilized and fine-tuned for the specific work to evaluate more accuracy and State-of-the-Art results.

Keywords- Multi-Lingual, Boundless Linguistic, Indian Languages, Natural Language Processing, BERT, Fine Tuning, State-of-the-Art, Supervised Algorithms.

1. INTRODUCTION:

India is a country with great linguistics and huge diversity which makes India a multilingual country. There are 456 languages spoken in India and 22 official [5] languages which makes it the richest language pool in the World. Marathi is an Indo-European language that is prominently spoken in Maharashtra. Marathi is the official language of the state of Maharashtra. But little work has been done in terms of processing, analyzing, and generating results as compared to the English language.

The main reason is corpus, embedding models, and pretrained models, libraries are not previously available. As a result, finding state-of-the-art results is difficult and the Marathi data which is available in huge quantities remains unprocessed and not utilized for any research purpose. For processing unstructured data which people are posting over the Internet in terms of their words, thoughts and views Such data needs to get processed to retrieve meaningfulness.

Natural Language Processing is a subfield of artificial intelligence (NLP)[2]. With NLP machine understands unstructured data and retrieve meaningful information from it. NLP methods are treasured for sentiment analysis. NLP is also known as computational linguistics.

Data collected for this research purpose is extracted from India's Maharashtrian cookery show 'Madhurasrecipe' and 'Ashwini's Kitchen Recipes' which never had traveled before in any research paper. For performing this task machine learning and deep learning approaches are used. State-of-the-art sentiments are generated from this data by categorizing text into three sentiments namely positive, negative, mixed positive-negative and neutral sentiments. Before the application of algorithms, Marathi Text is preprocessed and after this phase, the text is given as input to the multinomial Naïve Bayes algorithm.

2. PREVIOUS WORK:

Machine Learning[3] Supervised algorithms are used in many research papers to perform sentiment analysis and text classification of Indian languages. Deep Learning is also explored in some of the research papers.

Muhammad Abbas, et al. (2019) have proposed a Multinomial Naïve Bayes(MNB) classifier to detect sentiment analysis of movie reviews. The major focus is given on the number of occurrences of each word irrespective of word order which helps in identifying sentiments. TF-IDF method has been used in the research. Authors achieve a

significant result with MNB as compared to Naïve Bayes Classifier[1].

Sujata Deshmukh, et al. (2017) proposed a system for identifying hidden sentiments from Marathi language text. In the research corpus-based approach has been proposed. Marathi up-to-date corpus is formulated with the individual polarity of words similar to WordNet which is considered an English Corpus. Further, an algorithm is created to find the cumulative polarity of the text and to identify sentiments such as positive, negative, and neutral[10].

Snehal Pawar, et al. (2017) performs sentiment analysis using a Lexicon-based technique that is text containing positive and negative words previously defined. The algorithms applied are the SVM algorithm, Naïve Bayes, and Maximum Entropy Classifier. The author concludes that SVM gives better accuracy as compared to the other two algorithms[4].

Mohammed Ansari, et al. (2016) explore a simple approach, to finding sentiments from a mixed script written in English but the text is transliterated into Hindi and Marathi words (Hindi and Marathi words are written with the help of English characters). A language identification algorithm has been developed. To predict the sentiment after processing raw input using Natural Language Processing techniques, SVM and Random Forests classifier is implemented on the dataset. The author concludes that accuracy reaches 95% [6].

Mazhar Ali, et al. (2017) develops an SVM and K-model and implemented them on the dataset to perform sentiment analysis on Sindhi Text. Corpus is normalized and analyzed using Document Term Matrix (DTM) and Term Frequency-Inverse Document Frequency (TF-IDF). Text Structurization of Sindhi text constructed based on five questions using five Ws Who, What, Where, Why, and When. Further sentiment analysis is performed using Sindhi lexicons with Part of Speech tagging. Researchers stated that the Precision, Recall, F-score, and accuracy of the supervised model has given good results on the Sindhi corpus dataset[5].

Sonali Shah, et al. (2020) performed sentiment analysis on mixed text Marathi + English (Marglish) using parametric and non-parametric models. Logistic Regression (LR), Decision Tree (DCT), Bernoulli Naïve Bayes (BNB), Multilayer perceptron (MLP), etc. Out of the algorithms applied MLP and Bernoulli Naïve Bayes (BNB) give the best result according to the Authors. The results are confirmed with 10-fold cross-validation and statistical testing[9].

Atharva Kulkarni, et al. (2021) Marathi Sentiment Analysis Dataset L3CubeMahaSent which was made publicly available for researchers in the same domain. Tweets extracted consist of positive, negative, and neutral classes. Baseline classification is presented using deep learning

techniques CNN, LSTM, ULMFiT, and BERT-based deep learning models. The authors further reported that the best accuracy is constructed with CNN and IndicBERT with Indic fastText word embeddings[12].

3. CORPUS COLLECTION AND PROCESSING:

Data is collected from the Indian Marathi recipe youtube channel by extracting tweets using Beautiful Soup (BS4) which is a python data extraction library. Data collection consists of 9000 tweets out of which extracted features are 12 in number. Some English words which were used in tweets are transliterated to Marathi and then word embeddings have been figured out.

Initially, data is processed for Supervised learning by removing special characters, emoticons, numbers, etc. And generating word embeddings using iNLTK. This library is used for processing Marathi text. Word vectors generated by this process are given as the input to algorithms such as Multinomial Naïve Bayes. To create a classifier model using BERT, Unstructured text directly give as the input without processing it.

4. SYSTEM CONSTRUCTION

Supervised algorithm Naïve Bayes [3] is used in many research papers which perform binary classification either 0 or 1 that is Negative or positive. But in this research, some data values are mixed i.e., positive, and negative so it is required to calculate results based on such data. So Supervised Multinomial Naïve Bayes algorithm which considers feature vectors of the text is applied to the extracted data set.

BERT, which stands for Bidirectional Encoder Representations from Transformers comes under the unsupervised language representation model [12]. In this research four, BERT multilingual models are used to specially process Indian Languages such as Marathi. These Pretrained models are fine-tuned on the dataset to retrieve state-of-the-art results. The reason for fine-tuning is, Pretrained models are trained on specific data such as Wikipedia text. Using it directly for research-specific purposes may not give the accurate or required result. So fine-tuning the models gives better results in terms of accuracy.

This research paper presents supervised algorithms namely the Multinomial Naïve Bayes model and unsupervised techniques namely Bidirectional Encoder Representations from Transformers (BERT) models such as XLM-roberta-base, bert-base-multilingual-cased, distilbert-base-multilingual-cased and bert-base-multilingual-uncased. Multilingual models are used for sentiment analysis of Indian Marathi recipes collected data set. Pretrained BERT models

mentioned above are fine tuned to generate State of the Art results.

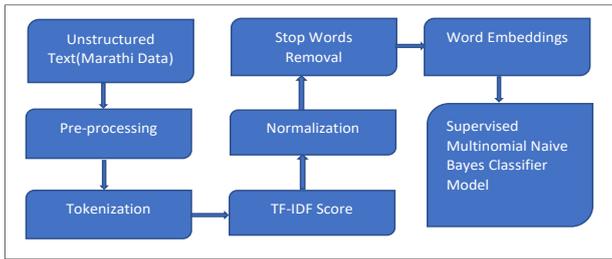


Figure 1: Sentiment Analysis using Modified Multinomial Naïve Bayes algorithm

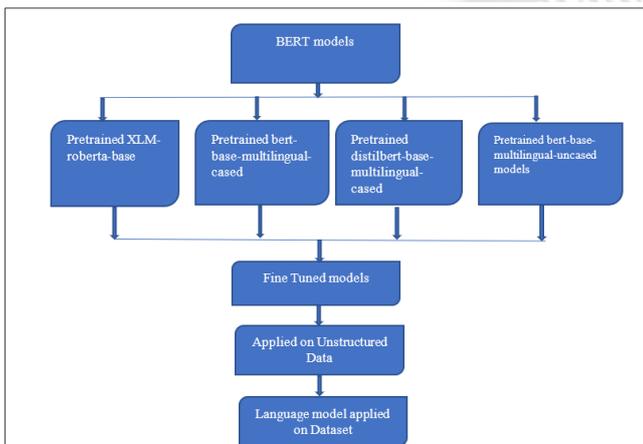


Fig 2: Sentiment Analysis using Fine-tuning BERT models

To perform sentiment analysis using a supervised multinomial Naïve Bayes classifier as mentioned in figure 1, unstructured data cannot directly provide to machine learning algorithms. Data in terms of text needed to be pre-processed first.

Pre-processing techniques executed on Marathi text are:

4.1. UNSTRUCTURED DATA -MARATHI TEXT:

Data is collected from the YouTube cookery channel which is an Indian YouTube channel specifically views posted by Maharashtrian people. This text consists of plain text, some English words, and some emoticons. Emoticons are removed from the text. Finding sentiments from Emoticons is the future scope of this research.

4.2. PRE-PROCESSING PHASE:

For supervised learning original raw text cannot be directly Provided to the machine, it needs pre-processing. Initially, raw text is processed to remove special characters and numbers like !, #, \$, %, ^, &, *, (,), ;, :, ', ', ', /, ?, etc., and numbers also which do not contribute any meaning to the sentence. Pre-processing is performed to make the text ready for Tokenization.

4.3. TOKENIZATION

Marathi Text is tokenized using the tokenizer Natural Language Processing Toolkit(nltk). Tokenizer will generate every sentence into a smaller unit of text. Tokens can be words, phrases, and n-grams that are extracted from the sentence. Execution for the sentence

data = "खरंच ताई तुम्ही लाजवाब आहे मस्त कस काय सुचतं नवीन नवीन पदार्थ बनवायला."

tokens will be

```
In [11]: from nltk.tokenize import sent_tokenize, word_tokenize
from nltk.corpus import stopwords
data = "खरंच ताई तुम्ही लाजवाब आहे मस्त कस काय सुचतं नवीन नवीन पदार्थ बनवायला "
stopwords = set(stopwords.words('marathi.txt')) # a set of Marathi stopwords
words1 = word_tokenize(data)

print(words1)

['खरंच', 'ताई', 'तुम्ही', 'लाजवाब', 'आहे', 'मस्त', 'कस', 'काय', 'सुचतं', 'नवीन', 'नवीन', 'पदार्थ', 'बनवायला']
```

Fig.3: Generating Tokens from Text

After tokenization stop words such "या", "आणि", "व", "यानी", "हे", "तर", "ते", "काही", "अशी", "असलेल्या" etc., are removed from the text. These stop words do not add any meaning to the text.

4.4. NORMALIZATION

After the stop words removal process, a stemming algorithm is not readily available and the Marathi stemmer algorithm model is created to remove suffix and prefix words that create noise and affects the expected result. This process is used to find the root word by which context meaning is generated.

Example sentence= "तुम्ही सांगितलेली शेगाव कचोरी खूपच छान झाली आणि पराठ्याची बनवण्याची सुद्धा कृती सांगणे."

After stemming results are

[' तुम्ही सांगितलेली शेगाव कचोरी खूप छान झाली आणि पराठा बनवणे सुद्धा कृती सांगा]

'उपवासाचे', 'पराठ्याचे' is normalized to root word ['उपवास', 'पराठा'] respectively.

4.5. TF-IDF FEATURES EXTRACTION:

Term Frequency — Inverse Document Frequency [4] is used for information retrieval or feature extraction. This process shows the importance of the word in a sentence.

It is calculated by using:

$$TF = \frac{\text{Occurrences of marathi word in a sentence}}{\text{Total count of words in a sentence}} \quad (i)$$

Example sentence= "तुम्ही सांगितलेली शेगाव कचोरी खूपच छान झाली .शेगाव कचोरी पराठ्याची बनवण्याची सुद्धा कृती सांगणे."

The values are further smoothened using log:

$$IDF = \log \left(\frac{\text{Total count of sentences}}{\text{frequency of word containing in each sentence}} \right) \quad (ii)$$

Features extracted according to TF-IDF score are: ['शेगाव', 'कचोरी']

4.6. WORD EMBEDDINGS:

Word embeddings are used to represent every single word that is analysed during the process of sentiment analysis. For English word embedding models are easily available such as word2vec, fasttext word embeddings, Glove word embeddings.

In this research, Natural Language Toolkit for Indic Languages (inltk) word embeddings are used and trained on Marathi Wikipedia text.

Word embeddings for Marathi text = “माझा शेगाव कचोरी बनवण्याचा खूप मोठा बिजनेस आहे.” are

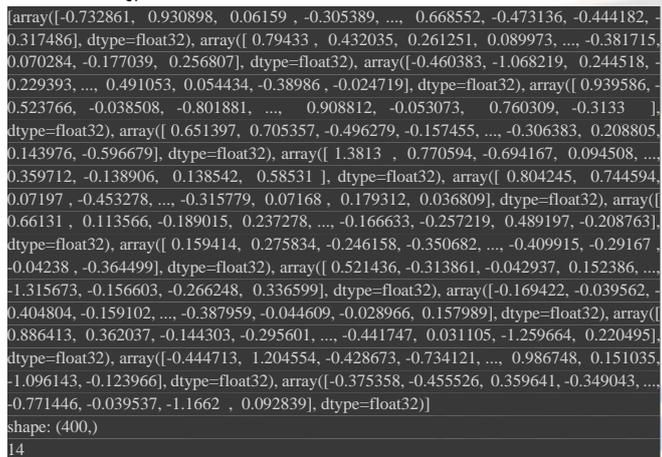


Fig 4: Word Embeddings for Marathi Text

5. ALGORITHMIC STRUCTURE

The first approach is to perform Multinomial Naïve Bayes (MNB) we need previously tagged sentences which is a supervised learning method. Sentences are classified into positive, negative, neutral, and mixed Targets.

Initially, sentence S is processed with Natural Language Processing techniques and converted into feature vector or embedding values E. Every feature e_i is the count with i^{th} term from S appeared or occurred in the corpus C, that has previously assigned to a tag T.

To compute the target T_k of sentence S:

Feature vectors are calculated from S denoted by W_{embv} . Each vector W_{embvi} .

The training dataset is previously tagged T_k using:

$$Pr(Wembv | Tk) = \frac{(\sum_{i=1}^n Wembvi)!}{\prod_{i=1}^n Wembvi} * \prod_{i=1}^n p_{ki}^{Wembvi} \quad (iii)$$

It is required to find the probability of tags Tk. It is calculated by the multinomial log space model:

$$\log Pr(Tk | Wembv) \propto \log p(Tk) + \sum_{i=1}^n Wembvi * \log p(Wembvi|Tk) \quad (iv)$$

From the above equation given the embedding vector Embvi of every word tag Tk is calculated.

Multinomial Naïve Bayes algorithm for sentiment analysis:

1. Sentence S is broken into to n terms using n grams which is unigram, bigram, or trigram.
2. For every kth tag T_k where $k=1$ to n perform:

3. Compute vector Wembvi for Wembv which is from tag Tk.
4. Find out prior probability $p(Tk)$ which occurred in a document from tag Tk.
5. Calculate posterior probability $Pr(Tk | Wembv)$ by adding prior $Pr(Tk)$ to the sum of each term Wembvi, given Tk:

$$Pr(Tk | Wembv) = \log p(Tk) + \sum_{i=1}^n Wembvi * \log p(Wembvi|Tk) \quad (iii)$$

The Second approach is Bidirectional Encoder Representations from Transformers [3] which is used to build Marathi classification model. BERT is an unsupervised algorithm. It considers the left and right context of the sentence as BERT can read sentences in both directions. Pretrained BERT is used and just one add on layer is created to predict state of the art sentiment for the dataset created. A single word may change the meaning of the sentence which is required to identify with the help of context that BERT does. It is not a sequential model, self-attention mechanism. BERT considers the effect of other words on sentence.

BERT multilingual base model [14] which can be used for Indian Languages such as Marathi such as XLM-roberta-base, BERT multilingual base model (cased), distilbert-base-multilingual-cased, bert-base-multilingual-uncased[15] are used in this research. Further these models are pretrained on Wikipedia text which consist of cleaned articles of Marathi language and other supported languages also are applied is fine tuned for sentiment analysis.

BERT is a transformer model for natural language processing. BERT identifies embeddings by considering text in both direction that why the name Bi-directional is given. Fine tuning BERT for sentiment analysis of Marathi Language gives remarkable performance. It is an unsupervised leaning method which do not require tagging to the text. BERT is largely deep bidirectional process. It takes into account context meaning of the statement unlike word2vec and glove for embedding generation. Executed code is

example Text=तुम्ही सांगितलेली शेगाव कचोरी छान झाली माझा मुंबईत शेगाव कचोरी बनवण्याचा खूप मोठा बिजनेस आहे.

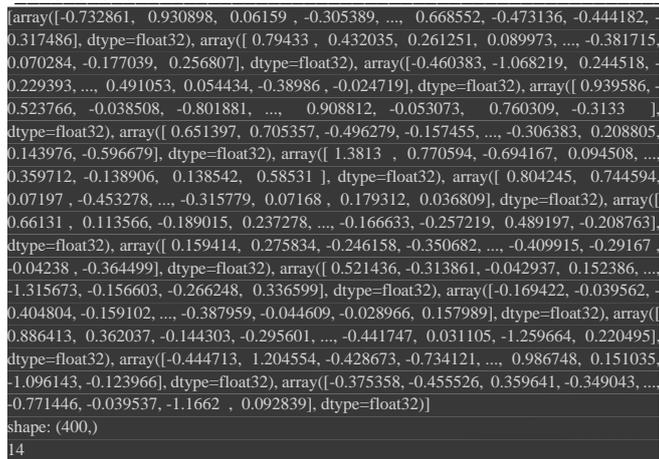


Fig 5: Context Based word embeddings generated by BERT

6. EXPERIMENTAL RESULTS:

Individual sentence polarity is given in the table.

User views.(Indian Marathi Text)	Supervised Learning Algorithm	Deep Learning Transformer				Sentiment Type
		MNB	BERT-XLM-Roberta-base	BERT multilingual base model (cased)	BERT distilbert-base-multilingual-cased	
वा फारच छान रेसिपी मी नक्की करते	61.2%	67.22%	68.32	68.35	69.01	positive
खरेच ताई तुम्ही लाजवाब आहे मस्त कस काय सुचतं नवीन नवीन पदार्थ बनवायला	58.70%	67.25	68.26	69.21	69.11	positive
हूनो टाकण्या ऐवजी दुसरा पर्याय आहे का	23.4%	36.33	37.32	38.45	38.47	negative
आम्हाला वाटत होता उपवास म्हणजे पोट रिकामे ठेवण, पोटाला आराम देणे	24.6%	35.23	37.37	36.44	37.55%	positive
पदार्थ छान आहे पण बनवायला जास्त वेळ आणि एनर्जी लागते रोजच्या थाळी बनवण्यापेक्षा	12.9%	32.31	32.41	33.33	34.44	mixed
चूक आहे हे त्यापेक्षा जेवण करा कि पोट भर उपवास बदनाम का करताय	10.3%	45.32	46.44	52.8%	52.22	negative

Table 1: Sentence Polarity

7. MODEL PERFORMANCE ON DATASET

Supervised and Transformer	Model Performance
Multinomial Naïve Bayes	57.33
BERT-XLM-roberta-base	65.44
bert-base-multilingual-cased	66.77
BERT-distilbert-base-multilingual-cased	66.89
bert-base-multilingual-uncased	68.90

Table 2: Model Performance

8. CONCLUSION and FUTURE WORK

Supervised and Unsupervised classifiers play a vital role in calculating and predicting Sentiments from Indian Languages. In this research when executing Sentiment analysis using BERT corpus size does not matter so we can

provide large paragraphs. But in many types of research accuracy is affected due to the text size. So out of the supervised learning classifier Multinomial Naïve Bayes and Unsupervised Language model, Bidirectional Encoder Representations from Transformers (BERT) produces state-of-the-art results. The accuracy generated by the BERT language model bert-base-multilingual-uncased is 68.90 which is good compared to the supervised model. Emoticons or emojis present in the text are not considered in this Research work, it can be further explored for Emoticons present in the text.

REFERENCES

- [1] M. Abbas, K. Memon, A. Jamali, S. Memon, and A. Ahmed, "Multinomial Naive Bayes Classification Model for Sentiment Analysis", IICSNS International Journal of Computer Science and Network Security, VOL.19 No.3, March 2019
- [2] K.Pal, "Automatic Multiclass Document Classification of Hindi Poems using Machine Learning Techniques", IEEE International Conference for Emerging Technology (INCET) Belgaum, India(2020)
- [3] A. Alsanad, "An Improved Arabic Sentiment Analysis Approach using Optimized Multinomial Naïve Bayes Classifier", International Journal of Advanced Computer Science and Applications, Vol. 13, No. 8, 2022
- [4] S. Pawar, and S. Mali, "International Journal on Recent and Innovation Trends in Computing and Communication" Volume: 5 Issue: 8
- [5] M. Ali and A. Wagan, "Sentiment Summarization and Analysis of Sindhi Text", (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 8, No. 10, 2017
- [6] Mohammed Arshad Ansari, and Sharvari Govilkar, "Sentiment Analysis of Transliterated Hindi and Marathi Script", 144 Sixth International Conference on Computational Intelligence and Information Technology – CIIT 2016
- [7] P. Sharma, and Teng-Sheng Moh, "Prediction of Indian Election Using Sentiment Analysis on Hindi Twitter", 2016 IEEE International Conference on Big Data (Big Data).
- [8] Kerstin Denecke, "Using SentiWordNet for Multilingual Sentiment Analysis", 2008 IEEE 24th International Conference on Data Engineering Workshop, 7-12 April 2008.
- [9] S. Shah, A. Kaushik, S. Sharma, and S.Sharma, "Opinion-Mining on Marglish and Devanagari Comments of YouTube Cookery Channels Using Parametric and Non-Parametric Learning Models", Big Data Cogn. Comput. 2020.
- [10] S. Deshmukh, N. Patil, S. Rotiwar, and J. Nunes, "Sentiment Analysis of Marathi Language", international journal of research publication in engineering and technology [IJRPET] ISSN 2454-7875, vol 3 ,issue 6, jun-2017
- [11] A. Sahani, K. Sarang, S.Umredkar, and M Patil, " Automatic Text Categorization of Marathi Language Documents", (IJCSIT) International Journal of Computer Science and

Information Technologies, Vol. 7 (5) , 2016, 2297-2301, ISSN:0975-9646

- [12] A. Kulkarni, M. Mandhane , M. Likhitar, G. Kshirsagar, and R. Joshi, "L3CubeMahaSent: A Marathi Tweetbased Sentiment Analysis Dataset", Proceedings of the 11th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, pages 213–220 April 19, 2021
- [13] C. Kariya, P. Khodke, "Twitter Sentiment Analysis", 2020 International Conference for Emerging Technology (INCET) Belgaum, India. Jun 5-7.
- [14] R. Naukarkar , Dr. A. Thakare, "A Design on Recognition of Sentiment Analysis of Marathi Tweets using Natural Language Processing", International Journal of Scientific Research in Science and Technology Print ISSN: 2395-6011 | Online ISSN: 2395-602X (www.ijrst.com) doi : <https://doi.org/10.32628/IJSRST>

