

Nepali News Classification using Naive Bayes: A Data-Driven Approach

KHADKA PILOT^{1*} AND POUDEL KSHITIZ^{2*}

¹Institute of Engineering, Thapathali Campus, Thapathali, Nepal (e-mail: pilot.076bct025@tcioe.edu.np)

²Institute of Engineering, Thapathali Campus, Thapathali, Nepal (e-mail: kshitizpoudel18@gmail.com)

Corresponding author: Khadka Pilot (e-mail: pilot.076bct025@tcioe.edu.np).

* Authors contributed equally

ABSTRACT The Nepali news landscape has witnessed an exponential growth in digital content, demanding mechanisms for news classification to help information retrieval and organization. This study explores the application of the Naive Bayes algorithm for the classification of Nepali news articles. The proposed approach leverages the simplicity and effectiveness of Naive Bayes, which assumes feature independence to achieve high classification accuracy while minimizing computational overhead. A comprehensive dataset of labeled Nepali news articles across diverse categories is utilized for model training and evaluation. The performance of the Naive Bayes classifier is assessed through various metrics, including precision and recall. The experimental results demonstrate promising outcomes, indicating the viability of the Naive Bayes approach for Nepali news classification tasks.

INDEX TERMS Naive Bayes, News classification

I. INTRODUCTION

OUR ancestors were motivated to find stable sources of food, water, and shelter to survive and protect their families. This deep need for certainty is true today as we work to address our basic needs and provide for ourselves and loved ones in a highly unpredictable world. The reality of our everyday existence is that things change every day, even when we wish that they would stay the same, and nothing is certain. Random events can and do happen when we're unprepared to deal with the fallout. [1]

This inherent uncertainty is not exclusive to our daily lives but can also be observed in the field of data mining. [2] Several reasons contribute to uncertainty in this domain:

- Inherent properties of the system being modeled: Take, for instance, quantum mechanics, where the behavior of subatomic particles is probabilistic.
- Incomplete observations: Weather forecasting serves as an example, where numerous variables interact, making it difficult to accurately include all variables.
- Incomplete modeling: While creating models, we often need to discard certain information. Such omissions can introduce uncertainty in our predictions, as seen in discretization of values.

Probability serves as a tool for quantifying uncertainty. One widely used technique that deals with uncertainty is the Naive Bayes classifier. This method relies on probabilities

and operates under the assumption of conditional independence between features to make predictions.

II. METHODOLOGY

A. THEORY

Bayes classifiers, based on Bayes theorem, are statistical classifiers that predict class membership probabilities.

The Naive Bayesian classifier assumes that the effect of an attribute value on a given class is independent of the values of other attributes. This assumption is called conditional independence.

Let $\mathbf{X} = (x_1, x_2, \dots, x_n)$ be a data tuple. In Bayesian terms, \mathbf{X} is considered as "evidence." Let H be the hypothesis that \mathbf{X} belongs to a specified class C . For classification, we want to determine $P(H|\mathbf{X})$, which represents the probability that hypothesis H holds true given tuple \mathbf{X} . [3]

The terms used in Bayes' theorem are as follows:

$P(H|\mathbf{X})$ is the posterior probability,

$P(H)$ is the prior probability,

$P(\mathbf{X}|H)$ is the likelihood,

$P(\mathbf{X})$ is the prior probability of \mathbf{X} .

Bayes' theorem provides a way of calculating the posterior probability $P(H|\mathbf{X})$ from $P(\mathbf{X}|H)$, $P(\mathbf{X})$, and $P(H)$, accord-

ing to the following equation:

$$P(H|X) = \frac{P(X|H) \cdot P(H)}{P(X)} \quad (1)$$

The numerator is focused on, as the denominator is constant for all classes. C_1, C_2, \dots, C_m . So, only $P(X|C_i) \cdot P(C_i)$ needs to be maximized. However, given the large number of attributes, it would be computationally expensive to calculate $P(X|C_i)$ directly. To reduce computation, the naive assumption is that of class conditional independence, i.e., there is no dependence relation among the attributes. This allows us to represent $P(X|C_i)$ as the product of the individual attribute probabilities:

$$P(X|C_i) = \prod_{k=1}^n P(x_k|C_i) \quad (2)$$

When we break text into words/features, multiplication of all likelihoods results in almost zero probabilities. To solve this we apply log which converts sum to product. Unlike multiplication logarithm is monotonously increasing function so the probabilities go on increasing.

$$\sum_{i=1}^n \log(p(x_i|C_k)) \quad (3)$$

B. TRAINING NAIVE BAYESIAN CLASSIFIER

For the class prior $P(C)$, the percentage of documents in training set that belong to each class C is determined. [4] Let N_c be the number of documents in our training data with class C , and N_{doc} be the total number of documents. Then, the class prior is calculated as follows:

$$P(C) = \frac{N_c}{N_{\text{doc}}} \quad (4)$$

To find $P(X|C)$, assumption is made that a feature is just the existence of a word w_i in the document's bag of words. Therefore, $P(w_i|c)$ can be computed as the fraction of times the word w_i appears among all words in all documents of topic C . Let V represent the vocabulary, which consists of the union of all word types in all classes, not just the words in one class C . The likelihood term is then calculated as follows:

$$P(w_i|C) = \frac{\text{count}(w_i, C)}{\sum_{w \in V} \text{count}(w, C)} \quad (5)$$

1) Zero frequency problem

Naive bayes faces zero frequency problem where it assigns zero probability to a word in test data set if it wasn't available in training dataset. To solve this smoothing is required.

C. ALGORITHM

D. SYSTEM BLOCK DIAGRAM

Appendix

Algorithm 1 Train Naive Bayes Classifier

```

1: function TrainNaiveBayes( $D, C$ )
2:    $N_{\text{doc}}$  = number of documents in  $D$ 
3:    $N_c$  = number of documents from  $D$  in class  $C$ 
4:    $\text{logprior}[c] \leftarrow \log \frac{N_c}{N_{\text{doc}}}$ 
5:    $V \leftarrow$  vocabulary of  $D$ 
6:    $\text{bigdoc}[c] \leftarrow \text{append}(d)$  for  $d \in D$  with class  $c$ 
7:   for each word  $w$  in  $V$  do
8:      $\text{count}(w, c) \leftarrow$  number of occurrences of  $w$  in  $\text{bigdoc}[c]$ 
9:      $\text{loglikelihood}[w, c] \leftarrow \log \left( \frac{\text{count}(w, c) + 1}{\sum_{w' \in V} (\text{count}(w', c) + 1)} \right)$ 
10:  end for
11:  return  $\text{logprior}, \text{loglikelihood}, V$ 
12: end function

```

Algorithm 2 Test Naive Bayes Classifier

```

1: function TestNaiveBayes( $\text{testdoc}, \text{logprior}, \text{loglikelihood}, C, V$ )
2:    $\text{sum} \leftarrow \{\}$ 
3:   for  $c \in C$  do
4:      $\text{sum}[c] \leftarrow \text{logprior}[c]$ 
5:     for  $i$  in 1 to length of  $\text{testdoc}$  do
6:        $\text{word} \leftarrow \text{testdoc}[i]$ 
7:       if  $\text{word} \in V$  then
8:          $\text{sum}[c] \leftarrow \text{sum}[c] + \text{loglikelihood}[\text{word}, c]$ 
9:       end if
10:    end for
11:  end for
12:  return  $\arg \max_{c \in C} \text{sum}[c]$ 
13: end function

```

E. INSTRUMENTATION

The following libraries were instrumental in completion of this project :

- **Pandas:** It is a library for data manipulation and analysis.
 - `df.drop`: Remove unnecessary columns.
 - `df.apply`: Remove punctuation, stopword from text.
- **Matplotlib:** It is a popular data visualization library for Python.
 - `plt.pie`: To visualize distribution of classes.
 - `plt.hist`: To visualize frequency of text length.
- **Scikit-learn:** Scikit-learn features algorithms like regression, classification, etc.
 - `train_test_split`: To split the dataset into training and test sets in the ratio 8:2.
 - `confusion_matrix`: To plot Confusion matrix after prediction.
 - `CountVectorizer, TfidfVectorizer`: To vectorize nepali text.
 - `GaussianNB, MultinomialNB, BernoulliNB`: To implement Naive Bayes for classification.
- **WordCloud:** It was used to visualize text data where the size of the word represents its frequency.

- **NLTK**: NLTK is a leading platform for building Python programs to work with human language data.
 - word_tokenize: To tokenize nepali text

III. WORKING PRINCIPLE

A. DATASET COLLECTION

The Nepali News classification dataset was gathered from diverse online news platforms. With a total of 7,470 instances, this comprehensive dataset is made up of three essential features: 'heading', 'paragraph' and 'label.' The 'heading' feature represents the concise title of the news, providing a quick glimpse into its content. On the other hand, the 'paragraph' feature delves deeper into the news, presenting a detailed description or body of the article. Finally, the 'label' attribute categorizes the news into one of three types:

- Business
- Sports
- Entertainment

B. DATA CLEANING

In this step, the data was examined for missing values, unexpected characters, and other anomalies. Out of the 7,470 instances, none of the values had null values. The title column was dropped. The reason behind this decision is that the 'paragraph' column already provides a more comprehensive and descriptive representation of the news content. Therefore, retaining the 'title' column would be redundant and offer no additional valuable information.

C. DATA PREPROCESSING

1) Remove punctuation

Various punctuation marks such as commas, colons, vertical bars, and others were thoughtfully removed from the text data. This punctuation removal process ensures that each text is treated equally during subsequent analysis and processing. Figure(2) shows the original sentence and Figure (3) shows the sentence after removing punctuation.

2) Stop word removal

Next, we address the removal of stop words, which are commonly occurring words used in sentence construction. In our dataset, we have identified several stop words such as shown in the figure(4). These words are deemed redundant for sentiment analysis and are therefore eliminated during the preprocessing stage.

3) Tokenization

Tokenization involves division of raw text into smaller unit known as tokens. These tokens can be individual words or entire sentences, depending on the chosen tokenizer. In our dataset, we opted for word tokenization. Figure (1) displays the original text before tokenization. After applying word tokenizer, Figure (2) exhibits the text transformed into a sequence of individual word tokens.

4) Vectorization

Vectorization or word embedding is a fundamental process in natural language processing that involves converting textual data into numerical vectors. By representing documents as numerical vectors, meaningful analytics and machine learning algorithms can be applied.

Term Frequency-Inverse Document Frequency (TF-IDF) is a NLP technique used to quantify the importance of a term in a specific document within a larger collection (corpus) of documents. TF represents the number of times a term appears in a particular document, making it document-specific.

- Raw Frequency: $tf(t)$ = Number of times term t occurs in a document.
- Term Frequency (Normalized):

$$(tf(t) = \frac{\text{Times term } t \text{ occurs in a document}}{\text{Total number of terms in the document}} \quad (6)$$

Inverse Document Frequency (IDF) is a measure of how common or rare a term is across the entire corpus of documents. It helps in distinguishing terms that are unique and significant within the corpus. The formula to calculate the IDF for a term t is:

$$idf(t) = 1 + \log_e \left(\frac{n}{df(t)} \right) \quad (7)$$

where: n is the total number of documents in the corpus. $df(t)$ is the number of documents in the corpus that contain the term t .

The TF-IDF value of a term in a document is the product of its TF and IDF. This combined value serves as a weight that indicates the relevance of the term in that specific document. A higher TF-IDF value implies that the term is more significant and relevant to that particular document compared to other terms.

D. CLASSIFICATION

The dataset is divided into two subsets in an 8:2 ratio, with 80% of the data used for training and the remaining 20% for testing. The training data is then utilized to train the Naive Bayes classifier. After the classifier has been trained, it is tested on the test set to evaluate its performance.

E. EVALUATION

Most real datasets aren't balanced. Even if the dataset is balanced, relying solely on accuracy may provide a misleading sense of model performance. In such cases, the confusion matrix becomes essential as it provides more comprehensive information.

Precision: Precision measures the proportion of cases that the model correctly flags as positive out of all the cases it predicted as positive. It helps us understand the accuracy of positive predictions.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \quad (8)$$

Recall: Recall (also known as sensitivity) measures the percentage of actual positive cases in the dataset that the model

correctly identified as positive. It helps us assess how well the model captures positive instances.

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} \quad (9)$$

F1-score: The F1 score is defined as the harmonic mean of precision (P) and recall (R) and is calculated using the following equation:

$$F1 = \frac{2 \times P \times R}{P + R} \quad (10)$$

By analyzing both precision and recall, we can gain a more nuanced understanding of the model's performance. High precision indicates that the model is making fewer false positive predictions, while high recall signifies that the model captures more positive cases correctly.

IV. RESULT AND ANALYSIS

A. DATASET ANALYSIS

The dataset consists of the following categories:

- Business: 2628 entries
- Sports: 2574 entries
- Entertainment: 2268 entries

Business: This category includes news related to business, economy, share market, and other economic activities.

Sports: This category covers various sports-related news, including competitions, tournaments, and other sporting events.

Entertainment: This category encompasses news related to movies and music, such as movie releases and music album launches.

B. RESULT ANALYSIS

Figure 7 shows distribution of training data in pie chart, and shows that each class has roughly equal number of instances, so our dataset is not skewed. Figure 8 and 9 shows the distribution of text length in dataset before and after removing stopwords. Stopwords are words that do not have much impact on the meaning of text. We see that the number of words have been reduced by removing stopwords.

Figure 10 shows wordcloud of all the words in our dataset. It shows the most frequent words in our dataset.

Figure 11, 12, 13 shows wordcloud for each class label. We can see that each class is characterized by few repeated words. Thus it is helpful for our naive bayes model to identify what sorts of feature fall into what class. Even if the features are assumed to be independent upon each other, we can classify text into corresponding classes by looking at individual words present.

Figure 14, 15 shows confusion matrix for classification. Both models have comparatively less accuracy and f1 score on entertainment class. The matrix also shows that the entertainment class is being confused for mostly sports class.

Table 1 shows the comparison of accuracy between different naive bayes and vectorization approaches. We can see that count vectorization performed significantly better than tf-idf

TABLE 1. Classification report with Multinomial NB classifier and counter vectorization

Class	Precision	Recall	F1-score	Support
sports	0.86	0.88	0.87	527
entertainment	0.79	0.80	0.80	428
business	0.90	0.87	0.88	539
Accuracy			0.85	1494
Macro avg	0.85	0.85	0.85	1494
Weighted avg	0.85	0.85	0.85	1494

TABLE 2. Classification report with Multinomial NB classifier and tfidf vectorization

Class	Precision	Recall	F1-score	Support
sports	0.67	0.80	0.73	532
entertainment	0.71	0.48	0.57	448
business	0.74	0.80	0.77	514
Accuracy			0.70	1494
Macro avg	0.71	0.69	0.69	1494
Weighted avg	0.71	0.70	0.70	1494

vectorization. This may be due to the fact that the texts are very short and also frequency of certain words matter very much to know the class of news. Also count vectorization provides straightforward and interpretable representation and doesn't introduce extra complexity which may confuse the model. We also observe that the Multinomial naive bayes is better since our representation consists of frequency of words. The gaussian naive bayes performs poor as our representation is sparse and also doesn't contain continuous attributes. The bernoulli naive bayes lies in between as it is equivalent to using bag-of-words approach which is not as good as counting the frequencies.

V. DISCUSSION

A. GAUSSIAN VS BERNOULLI VS MULTINOMIAL NAIVE BAYES

TABLE 3. Accuracy Comparison Chart

	Gaussian NB	Multinomial NB	Bernoulli NB
Count-vectorizer	79	85	84
TF-IDF vectorizer	54	70	64

Gaussian Naive Bayes is more suitable for dataset having continuous attribute values and where the distribution is Gaussian. Example weight, height of person etc. We fit Gaussian distribution to each feature and belonging to a particular class in which parameters (mean and variance) are determined by Maximum likelihood estimation. Then given a test instance, we calculate posterior probability of belonging to a particular class from all features. The class with highest probability is then assigned.

$$P(\mathbf{X}|C_k) = \prod_{i=1}^n p_{k_i}(x_i) (1 - p_{k_i})^{1-x_i} \quad (11)$$

Bernoulli Naive Bayes is suitable for data binary feature variables. For example if bag-of-words approach is used then

it is easier to use Bernoulli distribution to calculate probabilities.

$$p(\text{word}_i|\text{Class}_k) = \frac{\sum_{\text{training set}} \text{count}_{ik}}{\sum_{\text{training set}} \text{count}_k} \quad (12)$$

Unlike the bag of words approach we may want to count occurrence of each word/feature in text. **Multinomial Naive Bayes** is more suitable for dataset having frequencies of discrete data. Which makes it suitable for text classification with count vectorization.

B. SMOOTHING

As we are dealing with textual data, there is a great chance that we encounter unseen features in test instances. Not only that, the number of features is also not fixed. If new types of features occur, we can simply ignore them and see only those features encountered in training. However, this is an incorrect approach as it is equivalent to assigning a probability of 1. If not, then the probability becomes zero. To solve this, smoothing is required.

The Gaussian Naive Bayes doesn't require smoothing because it assumes all features are continuous. When it encounters a new feature value, it calculates the probability from a Gaussian distribution.

On the other hand, Categorical, Multinomial, and Bernoulli Naive Bayes require Laplace smoothing. Laplace smoothing calculates the probability as:

$$P_{\text{lap},k}(x|y) = \frac{c(x,y) + k}{c(y) + k|X|} \quad (13)$$

Where: x is the unseen word, y is the class label, $c(x,y)$ is the count of class ' y ' containing word ' x ', $c(y)$ is the total count of class ' y ', k is the smoothing factor, $|X|$ represents the number of features/dimensions in the data.

So, Laplace smoothing basically adds a new instance for each class containing the unseen feature value to the training data and calculates the probability for the required class.

C. POSSIBLE SOURCES OF ERRORS

- **Ambiguous or mislabeled data:** If the training dataset contains incorrectly labeled news articles or articles with ambiguous category assignments, the classifier might learn from these errors and make incorrect predictions during testing.
- **Vocabulary mismatch:** If the words present in the testing data are not present in the training data, the classifier may not know how to handle them, leading to misclassifications.
- **Out-of-vocabulary words:** New and previously unseen words (out-of-vocabulary words) in the testing data can cause issues if the model has not encountered them during training. These words will be ignored, leading to potential loss of information.
- **Dependencies between words:** Naive Bayes assumes that features (words) are conditionally independent, but in reality, words in a sentence often have dependencies.

Ignoring these dependencies may lead to errors and performance loss.

- **Uninformative features:** Some words might be present in multiple categories and not contribute much to distinguishing between classes. These uninformative features can confuse the classifier and lead to misclassifications.

VI. CONCLUSION

The Naive Bayes algorithm implementation on Nepali news classification demonstrated promising results. Through careful data preprocessing, including tokenization, stop word removal, and vectorization, the algorithm effectively transformed the textual data into a suitable format for analysis. The construction of a comprehensive vocabulary from the training dataset allowed for representation of the news articles as feature vectors.

Despite the simplicity of the Naive Bayes approach and its assumption of conditional independence between features, the classifier performed reasonably well in distinguishing between different categories of Nepali news. The model showcased its ability to handle large-scale datasets, making it suitable for real-world news classification applications.

However, there were certain challenges encountered during the analysis. Ambiguous or mislabeled data in the training set posed a risk of model learning from errors, leading to inaccurate predictions during testing. Additionally, the presence of out-of-vocabulary words in the testing data raised concerns about potential information loss and misclassifications. Furthermore, the assumption of independence between words in the Naive Bayes algorithm might not hold true for Nepali news articles, which could have influenced the model's performance.

To improve the classifier's accuracy, future work could focus on refining the data preprocessing techniques, addressing the issue of vocabulary mismatch, and exploring advanced algorithms that can capture word dependencies more effectively. Also, regular updates to model with new data would ensure that the classifier remains relevant and accurate in an evolving news landscape.

REFERENCES

- [1] Living with Uncertainty: How to Accept and Be More Comfortable with Unpredictability.// <https://news.stanford.edu/report/2020/11/04/living-with-uncertainty/>.
- [2] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *"Deep Learning."* MIT Press.
- [3] Han, J., Kamber, M., Pei, J. (2012). *Data Mining: Concepts and Techniques*, Third Edition. Morgan Kaufmann Publishers.
- [4] Daniel Jurafsky & James H. Martin. *Speech and Language Processing*.
- [5] M. Narasimha Murty, V. Susheela Devi. *Pattern Recognition: An Algorithmic Approach*.



PILOT KHADKA is a student at Institute of Engineering, Thapathali Campus. He is expected to graduate in Bachelor of Computer Engineering in 2024. During his time at Thapathali Campus, Pilot has actively engaged in various academic and extracurricular activities. He has participated in coding competitions, collaborated on software development projects, and demonstrated a keen interest in exploring new technologies.



KSHITIZ POUDEL is a student at Institute of Engineering, Thapathali Campus. He is expected to graduate in Bachelor of Computer Engineering in 2024. His relentless pursuit of knowledge and dedication to uplifting and empowering others make him an exceptional contributor and an invaluable asset to the academic community.

VII. APPENDIX

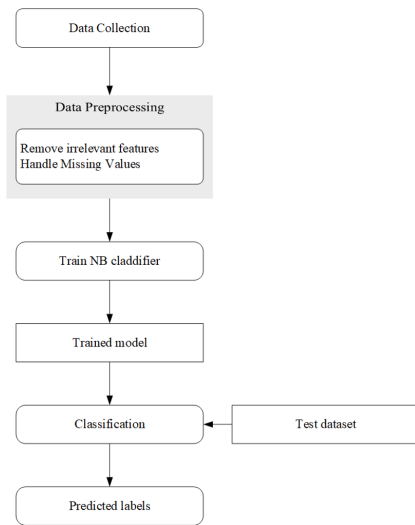


FIGURE 1. System Block Diagram

	headings	paras	label
0	कर्णालीका सुइना	उमेरले ७० कटेका पूर्णबहादुर विश्वकर्मा ७ वर्षीया बिरामी नातिनीको उपचार गर्न घरबाट हिँड्छन्। बाटोमा झमक्क सौँझ पर्छ। दायाँबायाँ छिटपुट घर छन् तर दलित भएकै कारण बास पाउँदैनन्। बल्लतल्ल गोठमा ओत लाग्न पाए पनि खानेकुरा दिइँदैन राति बिरामी नातिनी भोक लाग्यो भन्दै मुख बाउँछिन्।	entertainment
1	साकुराजस्तो प्रेमकथा	जापान भन्नेबित्तिकै पैयु साकुरा फूलको चर्चा हुन्छ।	entertainment
2	भद्रगोल नै ट्रेन्डिङमा	नेपाल टेलिभिजनबाट हरेक शुक्रबार सौँझ प्रसारण हुन्छ।	entertainment
3	फेरिए लोकभाका घर्नामा सामाजिक गीत	दोहोरी भन्नासाथ सोचिन्छ, यो केटा-केटीबीचको जुहुन्छ।	entertainment
4	भूकम्पले भत्केको सपना	च्यान्टे र मनमायाको स्थानीय परम्पराअनुसार विवाह हुन्छ।	entertainment

FIGURE 2. Nepali Sentiment Dataset

उमेरले ७० कटेका पूर्णबहादुर विश्वकर्मा ७ वर्षीया बिरामी नातिनीको उपचार गर्न घरबाट हिँड्छन्। बाटोमा झमक्क सौँझ पर्छ। दायाँबायाँ छिटपुट घर छन् तर दलित भएकै कारण बास पाउँदैनन्। बल्लतल्ल गोठमा ओत लाग्न पाए पनि खानेकुरा दिइँदैन। राति बिरामी नातिनी भोक लाग्यो भन्दै मुख बाउँछिन्।

FIGURE 3. Original Sentence

उमेरले ७० कटेका पूर्णबहादुर विश्वकर्मा ७ वर्षीया बिरामी नातिनीको उपचार गर्न घरबाट हिँड्छन्। बाटोमा झमक्क सौँझ पर्छ। दायाँबायाँ छिटपुट घर छन् तर दलित भएकै कारण बास पाउँदैनन्। बल्लतल्ल गोठमा ओत लाग्न पाए पनि खानेकुरा दिइँदैन राति बिरामी नातिनी भोक लाग्यो भन्दै मुख बाउँछिन्।

FIGURE 4. After removing punctuation

'छ', 'र', 'पति', 'छन्', 'लागि', 'भएको', 'गरेको'

FIGURE 5. Sample stopwords

उमेरले ७० कटेका पूर्णबहादुर विश्वकर्मा ७ वर्षीया बिरामी नातिनीको उपचार गर्न घरबाट हिँड्छन्। बाटोमा झमक्क सौँझ पर्छ। दायाँबायाँ छिटपुट घर छन् तर दलित भएकै कारण बास पाउँदैनन्। बल्लतल्ल गोठमा ओत लाग्न पाए पनि खानेकुरा दिइँदैन राति बिरामी नातिनी भोक लाग्यो भन्दै मुख बाउँछिन्।

FIGURE 6. After removing stopwords

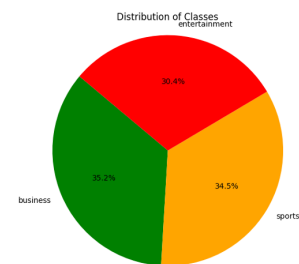


FIGURE 7. Distribution of classes

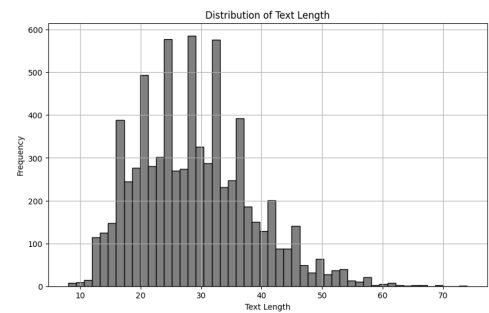


FIGURE 8. Text length histogram

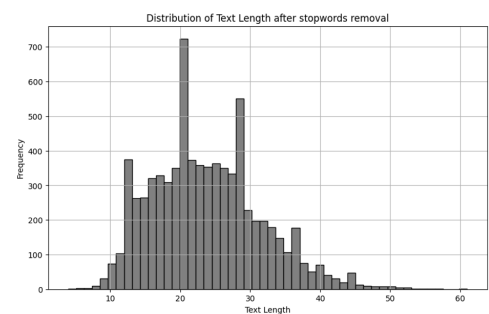


FIGURE 9. Text length histogram after stopwords removal



VIII. CODE

```

1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4
5  from google.colab import drive
6  drive.mount('/content/drive')
7
8  dataset_path = '/content/drive/MyDrive/
   Colab_Notebooks/Nepali_News_Classification.csv'
9
10 # Load the dataset into a pandas DataFrame
11 df = pd.read_csv(dataset_path)
12
13 df.head(10)
14
15 #check for null values
16 df.isnull()
17
18 df.info()
19
20 """##Visualizing the dataset"""
21
22 df.shape
23
24 class_counts = df['label'].value_counts()
25 print(class_counts)
26
27 #remove title column
28 df = df.drop(df.columns[0], axis=1)
29
30 df.head()
31
32 df.sample(frac=1).reset_index(drop=True)
33
34 class_counts = df['label'].value_counts()
35 print(class_counts)
36 # Create a pie chart
37 plt.figure(figsize=(6, 6))
38 plt.pie(class_counts, labels=class_counts.index,
   autopct='%1.1f%%', startangle=140, colors=['
   green', 'orange', 'red'])
39 plt.axis('equal')
40
41 # Add a title
42 plt.title('Distribution of Classes')
43
44 # Display the pie chart
45 plt.show()
46
47 df['paras'][0]
48
49 len(df['paras'][0])
50
51 #punctuation removal
52 import re
53 # Define a function to remove unnecessary
   punctuations from Nepali text
54 def remove_punctuation(text):
55     # Define the regular expression pattern to
   remove punctuations
56     punctuation_pattern = r'[ ?,:;\',.\(\)\n&
   !-]' # Add other punctuations
   as needed
57
58     # Use regular expression to remove
   punctuations
59     return re.sub(punctuation_pattern, '', text)
60
61 # Apply the function to remove punctuations from
   the 'Sentences' column
62 df['paras'] = df['paras'].apply(remove_punctuation

```

```

   )
63
64 print(df['paras'][0])
65
66 len(df['paras'][0])
67
68 #define a proper function to calculate words in a
   nepali text.
69 def get_length(text):
70     # Split the text into words based on spaces
71     words = text.split()
72
73     # Return the number of words in the text
74     return len(words)
75
76 get_length(df['paras'][0])
77
78 # Apply the function to calculate the text length
   for each instance
79 text_length = df['paras'].apply(get_length)
80
81 # Plot the text length distribution
82 plt.figure(figsize=(10, 6))
83 plt.hist(text_length, bins=50, color='gray',
   edgecolor='black')
84 plt.xlabel('Text Length')
85 plt.ylabel('Frequency')
86 plt.title('Distribution of Text Length')
87 plt.grid(True)
88 plt.show()
89
90 #stopwords removal
91
92 # Function to remove Nepali stopwords from a
   sentence
93 def remove_stopwords(sentence):
94     words = sentence.split()
95     filtered_sentence = [word for word in words if
   word.lower() not in nepali_stopwords]
96     return ' '.join(filtered_sentence)
97
98 # Create a new DataFrame with processed 'Sentences'
   (Nepali stopwords removed)
99 new_df = df.copy()
100 new_df['paras'] = new_df['paras'].apply(
   remove_stopwords)
101
102 new_df.head()
103
104 new_df['paras'][0]
105
106 # Apply the function to calculate the text length
   for each instance
107 text_length = new_df['paras'].apply(get_length)
108
109 # Plot the text length distribution
110 plt.figure(figsize=(10, 6))
111 plt.hist(text_length, bins=50, color='gray',
   edgecolor='black')
112 plt.xlabel('Text Length')
113 plt.ylabel('Frequency')
114 plt.title('Distribution of Text Length after
   stopwords removal')
115 plt.grid(True)
116 plt.show()
117
118 print(np.max(np.array(text_length)),',', np.min(np.
   array(text_length)))
119
120 """##Wordcloud"""
121
122 font_path="/content/drive/MyDrive/Colab_Notebooks/
   Mangal400.TTF"
123

```

```

124 from wordcloud import WordCloud
125 # Combine all the text data into a single string
126 text_data = ' '.join(new_df['paras'])
127 # Create a WordCloud object with the custom font
128 wordcloud = WordCloud(width=800, height=400,
129                        background_color='white', font_path=font_path,
130                        regexp=r"[\u0900-\u097F]+").generate(text_data,
131                        )
132 # Display the word cloud
133 plt.figure(figsize=(10, 6))
134 plt.imshow(wordcloud, interpolation='bilinear')
135 plt.axis('off')
136 plt.title('Word Cloud of Entire DataFrame')
137 plt.show()
138
139 #word cloud of preprocessed data
140 sentiments = new_df['label'].unique()
141
142 # Function to create and display a WordCloud for a
143 # specific sentiment category
144 def generate_wordcloud_for_sentiment(sentiment):
145     # Filter the DataFrame for the current
146     # sentiment
147     filtered_df = new_df[new_df['label'] ==
148                          sentiment]
149
150     # Combine all the text data for the current
151     # sentiment into a single string
152     textual_data = ' '.join(filtered_df['paras'])
153
154     # Create a WordCloud object with the custom
155     # font
156     wordcloud = WordCloud(width=800, height=400,
157                          background_color='white', font_path=font_path,
158                          regexp=r"[\u0900-\u097F]+").generate(
159         textual_data)
160
161     # Display the WordCloud for the current
162     # sentiment
163     plt.figure(figsize=(10, 6))
164     plt.imshow(wordcloud, interpolation='bilinear')
165     plt.axis('off')
166     plt.title(f'Word Cloud for Class: {sentiment}')
167     plt.show()
168
169 # Generate WordClouds for each sentiment category
170 for sentiment in sentiments:
171     generate_wordcloud_for_sentiment(sentiment)
172
173 """##Vectorizing the textual data through count-
174 vectorizer"""
175 new_df=new_df.sample(frac=1).reset_index(drop=True)
176
177 new_df.head()
178
179 from sklearn.feature_extraction.text import
180 CountVectorizer
181 # Create an instance of TfidfVectorizer
182 count_vectorizer= CountVectorizer(max_features
183 =300)#most repeated 300 words
184
185 # Fit and transform the text data using the
186 # vectorizer
187 X1=count_vectorizer.fit_transform(new_df['paras'])
188 .toarray()
189
190 print(X1[0],X1[0].shape)#6238 unique words
191
192 # Convert the NumPy ndarray to a pandas DataFrame
193 X1_df = pd.DataFrame(X1)
194
195 X1_df.head()
196
197 y=new_df['label']
198
199 print(y)
200
201 from sklearn.model_selection import
202 train_test_split
203 X_train1, X_test1, y_train1, y_test1 =
204     train_test_split(X1_df, y, test_size=0.2,
205                     random_state=42)
206
207 X_train1.shape
208
209 y_train1.shape
210
211 from sklearn.naive_bayes import GaussianNB,
212 MultinomialNB, BernoulliNB
213 modelGNB=GaussianNB()
214 modelMNB=MultinomialNB()
215 modelBNB=BernoulliNB()
216
217 modelGNB.fit(X_train1,y_train1)
218 modelMNB.fit(X_train1,y_train1)
219 modelBNB.fit(X_train1,y_train1)
220
221 y_pred1=modelGNB.predict(X_test1)
222 y_pred2=modelMNB.predict(X_test1)
223 y_pred3=modelBNB.predict(X_test1)
224
225 y_pred1
226
227 from sklearn.metrics import confusion_matrix,
228 classification_report
229 def get_report_and_cm(name):
230     cm = confusion_matrix(y_test1, name)
231
232     # Create a classification report
233     class_names = ['sports', 'entertainment',
234                   'business']
235     cr = classification_report(y_test1, name,
236                               target_names=class_names)
237
238     # Print the confusion matrix and classification
239     # report
240     print("Confusion Matrix:")
241     print(cm)
242
243     print("\nClassification Report:")
244     print(cr)
245
246 get_report_and_cm(y_pred1)
247 get_report_and_cm(y_pred2)
248 get_report_and_cm(y_pred3)
249
250 class_names = ['sports', 'entertainment', 'business']
251
252 import seaborn as sns
253 plt.figure(figsize=(10, 8))
254 sns.heatmap(confusion_matrix(y_test1, y_pred2),
255             annot=True, fmt='d', cmap='Blues', xticklabels
256             =class_names, yticklabels=class_names)
257 plt.xlabel('Predicted Labels')
258 plt.ylabel('True Labels')
259 plt.title('Confusion Matrix')
260 plt.show()
261
262 """##Vectorizing textual data through TF-IDF"""
263
264 new_df.head()
265
266 from sklearn.feature_extraction.text import
267 TfidfVectorizer

```

```

240 from nltk.tokenize import word_tokenize
241 from nltk.corpus import stopwords
242 import string
243 import nltk
244
245 nltk.download('punkt')
246 nltk.download('stopwords')
247
248 all_text = new_df['paras'].tolist()
249 # Join all the text elements into a single string
250 combined_text = ' '.join(all_text)
251
252 # Tokenize the combined text into words
253 words = word_tokenize(combined_text)
254
255 #remove stopwords
256 stop_words = set(stopwords.words('nepali'))
257 words = [word for word in words if word not in
258          stop_words]
259
260 # Create the vocabulary (set of unique words)
261 vocabulary = set(words)
262
263 print(vocabulary)
264
265 print(len(vocabulary))
266
267 def tokenize_text(text):
268     words = word_tokenize(text)
269     words = [word for word in words if word not in
270             stop_words]
271     return words
272
273 # Apply tokenization to the 'text' column in the
274 # DataFrame
275 new_df['tokenized_text'] = new_df['paras'].apply(
276     tokenize_text)
277
278 new_df.head()
279
280 #finally apply tf-idf vectorization
281 # Convert tokenized_text back to strings from
282 # lists of words
283 new_df['tokenized_text'] = new_df['tokenized_text']
284     .apply(lambda tokens: ' '.join(tokens))
285
286 # Initialize the TF-IDF vectorizer with your
287 # vocabulary
288 tfidf_vectorizer = TfidfVectorizer(vocabulary=
289     vocabulary, max_features=500)
290
291 # Fit and transform the tokenized_text to get the
292 # TF-IDF vectors
293 tfidf_vectors = tfidf_vectorizer.fit_transform(
294     new_df['tokenized_text'])
295
296 tfidf_df = pd.DataFrame(tfidf_vectors.toarray(),
297     columns=tfidf_vectorizer.get_feature_names_out(
298         ))
299
300 # Concatenate the original DataFrame with the TF-
301 # IDF DataFrame
302 new_df = pd.concat([new_df, tfidf_df], axis=1)
303
304 print()
305
306 new_df.head()
307
308 np.max(np.array(tfidf_df.head(1)))
309
310 X_train2, X_test2, y_train2, y_test2 =
311     train_test_split(tfidf_df, y, test_size=0.2,
312                     random_state=42)
313
314 modelGNB.fit(X_train2,y_train2)
315
316 modelMNB.fit(X_train2,y_train2)
317 modelBNB.fit(X_train2,y_train2)
318
319 y_pred21=modelGNB.predict(X_test2)
320 y_pred22=modelMNB.predict(X_test2)
321 y_pred23=modelBNB.predict(X_test2)
322
323 get_report_and_cm(y_pred21)
324 get_report_and_cm(y_pred22)
325 get_report_and_cm(y_pred23)
326
327 import seaborn as sns
328 plt.figure(figsize=(10, 8))
329 sns.heatmap(confusion_matrix(y_test2, y_pred22),
330     annot=True, fmt='d', cmap='Blues', xticklabels
331     =class_names, yticklabels=class_names)
332 plt.xlabel('Predicted Labels')
333 plt.ylabel('True Labels')
334 plt.title('Confusion Matrix')
335 plt.show()

```

...