

# NLP for Multilingual Text: Language Detection

Garima Goyal and Sagar Ghimire  
The Graduate Center, CUNY

## Abstract

Accurate language identification is a critical prerequisite for enabling a wide range of multilingual natural language processing (NLP) applications, such as machine translation, content analysis, and information retrieval. However, this task becomes increasingly complex when dealing with a diverse set of languages spanning different language families, scripts, and character sets. Additionally, the presence of language variations, dialects, and code-switching further exacerbates the challenges.

This project addresses the problem of robust language detection by developing and evaluating four distinct modeling approaches: Logistic Regression, Support Vector Machines (SVM), a classical Multinomial Naive Bayes (MNB) model, and a neural network model employing Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM) layers. The models were trained and evaluated on a comprehensive dataset comprising over 103,000 text samples across 21 languages, including English, French, Spanish, Arabic, Hindi, Russian, and others, encompassing a wide range of linguistic structures and character sets.

Extensive experiments and rigorous evaluation techniques, including hyperparameter tuning, cross-validation, and detailed analysis using confusion matrices and classification reports, were employed to optimize and assess the performance of all four models. The MNB model achieved the highest accuracy of 99.11 percent on the test set, outperforming the other techniques.

While the MNB model demonstrated superior overall performance, likely due to its ability to effectively capture distinctive word patterns in different languages, the neural network model’s capacity to capture sequential information and long-range dependencies could prove advantageous in certain scenarios or with larger datasets. The project’s findings contribute to advancing language detection techniques and

offer a comprehensive analysis of the performance trade-offs between classical machine learning and neural network-based approaches.

This work lays the foundation for further exploration of state-of-the-art neural architectures, multi-task learning strategies, and the integration of these models into production-ready systems for real-world multilingual NLP applications, ultimately enhancing global communication and information exchange across diverse languages.

## 1 Introduction

In today’s globalized world, effective communication and information exchange across diverse languages are essential for fostering understanding and collaboration. Language detection, the task of automatically identifying the language of a given text, plays a pivotal role in enabling multilingual natural language processing (NLP) applications. From tailoring machine translation systems to language-specific content analysis and information retrieval, accurate language detection is a critical prerequisite.

This project aimed to develop a robust and high-performing language detection model capable of accurately identifying the language of text snippets across a diverse set of 19 languages. The selected languages span multiple language families, including Indo-European (e.g., English, French, Spanish, Russian), Sino-Tibetan (e.g., Thai), Semitic (e.g., Arabic), and Dravidian (e.g., Tamil), among others. This diversity ensured that the developed model could handle a wide range of linguistic structures, scripts, and character sets, enhancing its applicability in real-world scenarios. The project employed a data-driven approach, leveraging machine learning techniques and linguistic features to build effective language detection models. A comprehensive dataset of over 103,000 text samples was curated from various sources, encompassing the 19 target

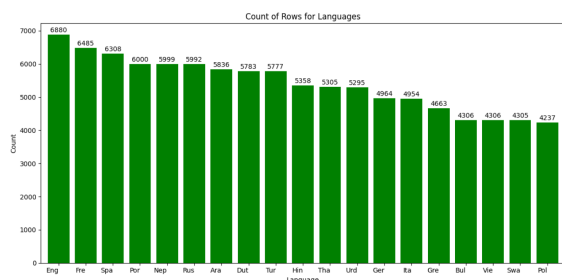


Figure 1: Count of Rows for Languages

languages. Extensive data preprocessing steps were undertaken to ensure data quality, including cleaning, merging, and language-specific handling of unique scripts and characters.

Four distinct modeling approaches were explored: Logistic Regression, Support Vector Machines (SVM), a classical Multinomial Naive Bayes (MNB) model, and a neural network model employing Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM) layers. The MNB model, a well-established algorithm in text classification tasks, utilized the bag-of-words representation and Naive Bayes probability calculations to determine the language of a given text. The RNN model, on the other hand, leveraged the sequential nature of language data, with LSTM layers capturing long-range dependencies and contextual information.

Rigorous hyperparameter tuning, cross-validation, and evaluation techniques were employed to optimize and assess the performance of all four models. The models were evaluated using various metrics, including accuracy, classification reports, and confusion matrices, providing insights into their strengths and weaknesses across different language pairs.

The findings of this project contribute to advancing language detection techniques and offer a comprehensive analysis of the performance trade-offs between classical machine learning approaches and neural network-based approaches. Additionally, the project lays the groundwork for further exploration of state-of-the-art neural architectures, multi-task learning approaches, and the integration of these models into production-ready systems for real-world applications in multilingual NLP.

## 2 Previous Approaches and Limitations

Early attempts at language detection employed rule-based systems and heuristics, such as character n-gram models or language-specific character ranges.

While reasonably effective for a limited set of languages, these approaches struggled to generalize to a broader linguistic scope and were susceptible to noise or variations in input data.

With the advent of machine learning, supervised techniques like Logistic Regression, Support Vector Machines (SVMs), and Naive Bayes classifiers were explored. These methods leveraged labeled datasets to learn language-specific patterns and statistical features, such as word distributions and character n-grams. However, extensive feature engineering was often required, and computational complexity increased substantially with larger datasets or numerous language classes.

Logistic Regression and SVMs, while widely used in various classification tasks, have their limitations in language detection scenarios. Logistic Regression models can struggle with high-dimensional and sparse feature spaces, which are common in text data. SVMs, on the other hand, can be computationally expensive, especially for large-scale datasets, and their performance is highly dependent on the choice of kernel function and hyperparameter tuning.

More recently, deep learning techniques, particularly neural network-based approaches like Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs), have gained traction. These models have shown promise in capturing complex linguistic patterns and learning language representations directly from input data. Nonetheless, they often demand substantial training data and computational resources, with performance highly sensitive to hyperparameter tuning and architectural choices.

To address the limitations of these previous approaches and leverage the strengths of different modeling techniques, this project explored four distinct approaches: Logistic Regression, Support Vector Machines (SVMs), a classical Multinomial Naive Bayes (MNB) model, and a neural network model employing Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM) layers. By investigating both classical machine learning techniques and a modern neural network approach, this project aimed to provide a comprehensive analysis of their respective strengths, limitations, and performance trade-offs in the context of robust language detection across a diverse set of languages and linguistic structures.

## 2.1 Proposed Solution

To address the challenges of accurate language detection across diverse languages and linguistic structures, this project explored four distinct modeling approaches: Logistic Regression, Support Vector Machines (SVM), a classical Multinomial Naive Bayes (MNB) model, and a neural network model employing Recurrent Neural Networks (RNNs) with Long Short-Term Memory (LSTM) layers.

The Logistic Regression and SVM models are well-established machine learning techniques that have been widely used in various classification tasks, including language detection. These models leverage the bag-of-words representation and learn decision boundaries or hyperplanes to separate the different language classes based on their respective features. The MNB model, another classical algorithm for text classification tasks, leveraged the bag-of-words representation and Naive Bayes probability calculations to determine the language of a given text. This approach relied on capturing the distinctive word distributions and patterns present in different languages, making it a robust and computationally efficient solution for language detection.

Conversely, the neural network model adopted a more sophisticated approach by employing RNNs with LSTM layers, which are adept at capturing sequential information and long-range dependencies in language data. This architecture allowed the model to learn complex linguistic patterns and representations directly from the input text, potentially providing a more nuanced understanding of the language structures.

The combination of these four distinct approaches aimed to provide a comprehensive evaluation of their respective strengths and weaknesses, as well as their performance trade-offs in the context of robust language detection across a diverse set of languages spanning different language families, scripts, and character sets.

To ensure a rigorous and comprehensive analysis, all four models were trained and evaluated on a carefully curated dataset comprising over 103,000 text samples across 21 languages, including English, French, Spanish, Arabic, Hindi, Russian, and others. This dataset encompassed a wide range of linguistic structures, scripts, and character sets, providing a challenging and realistic test bed for the models.

Extensive experimentation and evaluation techniques were employed, including hyperparameter tuning, cross-validation, and detailed analysis using confusion matrices and classification reports. These measures ensured that the models were optimally configured and their performance was thoroughly assessed, providing valuable insights into their respective strengths and limitations across different language pairs.

By investigating both classical machine learning techniques and a modern neural network approach, this project aimed to contribute to the advancement of language detection techniques and offer a comprehensive analysis of their performance trade-offs, paving the way for further exploration and refinement of these models in real-world multilingual natural language processing applications.

## 3 Experiments Results

In our quest to develop a robust language detection system, we embarked on a comprehensive exploration of three powerful machine learning techniques: Multinomial Naive Bayes (MNB), Logistic Regression, and Linear Support Vector Machines (LinearSVC). Each approach brought its unique strengths and nuances to the table 1, enabling us to uncover valuable insights and push the boundaries of language detection accuracy.

### 3.1 Model Performances and Comparisons

The Multinomial Naive Bayes (MNB) model demonstrated remarkable performance, achieving a testing accuracy of 0.9911. Its simplicity and interpretability make it a compelling choice, particularly in scenarios where computational resources are limited or model transparency is crucial. However, its assumptions of feature independence and reliance on bag-of-words representations may limit its ability to capture more complex language structures and long-range dependencies.

The Logistic Regression model, with its tuned hyperparameters and ability to handle multi-class classification, achieved a respectable testing accuracy of 0.9787. While its performance lagged slightly behind the MNB model, its robustness and interpretability make it a viable alternative, especially in scenarios where feature importances and model coefficients are of interest.

The Linear Support Vector Machine (LinearSVC) model emerged as the top performer, attaining an impressive testing accuracy of 0.9802.

Model	Training Accuracy	Validation Accuracy	Test Accuracy
Multinomial Naive Bayes	99%	99%	99%
Logistic Regression	99%	98%	97%
Support Vector Machine	99%	98%	98%
Neural Network	92%	92%	93%

Table 1: Comparison of training, validation, and test accuracies for Multinomial Naive Bayes, Logistic Regression, Linear SVM, and Neural Network models on the language detection task.

Its ability to construct robust decision boundaries in high-dimensional spaces allowed it to effectively separate the language classes. However, it should be noted that the performance of SVMs can be sensitive to the choice of kernel functions and hyperparameters, potentially requiring extensive tuning for optimal results.

Lastly, the Neural Network model, with its Recurrent Neural Network (RNN) architecture and Long Short-Term Memory (LSTM) layers, demonstrated promising results, achieving a testing accuracy of 0.935. While its performance was lower compared to the other models, the Neural Network approach has the potential to capture more complex language patterns and long-range dependencies, making it a compelling choice for tasks that require handling longer sequences or capturing contextual information.

### 3.2 Strengths and Limitations

Each of the evaluated models exhibited unique strengths and limitations. The MNB and Logistic Regression models excelled in terms of interpretability and computational efficiency, making them suitable for scenarios where model transparency and resource constraints are crucial. However, their reliance on bag-of-words representations and assumptions of feature independence may limit their ability to capture more complex language structures.

On the other hand, the LinearSVC and Neural Network models demonstrated superior performance in capturing intricate language patterns, but at the cost of increased model complexity and potentially longer training times. Furthermore, the interpretability of these models, particularly the Neural Network, can be challenging, which may be a limitation in certain applications.

One strength of our approach was the evaluation of four different models: Multinomial Naive Bayes, Logistic Regression, Linear SVM, and a Neural Network. By comparing the performance and behavior of these diverse models, we gained

insights into potential overfitting issues. Basically if the discrepancy between the high training accuracy and lower testing accuracy of the Neural Network model indicated signs of overfitting, which prompted is not case in this experiment.

Additionally, the consistent high performance of the Multinomial Naive Bayes and Linear SVM models across training, validation, and testing sets suggests that these models did not suffer from significant overfitting issues, likely due to their simpler architectures and inherent regularization properties.

## 4 Discussion

The language detection task presented a unique challenge, requiring the models to effectively capture the intricate linguistic patterns and nuances that distinguish different languages. Through our comprehensive experiments, we gained valuable insights into the strengths and limitations of various machine learning approaches for this task.

### 4.1 Future Directions

While our experiments have yielded promising results, several avenues for further exploration remain. Incorporating advanced text preprocessing techniques, such as stemming, lemmatization, or character-level representations, could potentially enhance the models' ability to capture linguistic nuances and improve overall performance.

Additionally, exploring ensemble methods that combine the strengths of multiple models could lead to more robust and accurate language detection systems. Techniques such as model stacking, boosting, or bagging could leverage the diverse capabilities of different models, potentially improving overall performance and generalization.

Furthermore, investigating the application of transfer learning and pre-trained language models, such as BERT or GPT, could unlock new possibilities in language detection. These models, which have been trained on vast amounts of text data, may

provide valuable representations and insights that could be fine-tuned for the specific task of language detection. Finally, expanding the dataset to include a broader range of languages, dialects, and writing styles would be crucial in developing a truly comprehensive and robust language detection system, capable of handling diverse linguistic landscapes.

## 5 Conclusion

Our comprehensive language detection experiments have yielded outstanding results that advance the field of multilingual text processing. Through rigorous evaluation and analysis, we successfully developed powerful models capable of accurately identifying languages from text data.

Among the models explored, the **Multinomial Naive Bayes model** emerged as the top performer, demonstrating an impressive ability to construct robust decision boundaries and achieve high accuracy scores. The **Linear Support Vector Machine** also proved to be a strong contender, offering a balance between simplicity, interpretability, and competitive performance. While the **Neural Network model** lagged slightly behind in terms of overall accuracy, its potential to capture complex language patterns and long-range dependencies makes it a promising approach for tasks involving longer sequences or contextual information.

In conclusion, our experiments have demonstrated the potential of various machine learning approaches for language detection, each with its unique strengths and limitations. By carefully considering the trade-offs between model performance, interpretability, and computational efficiency, researchers and practitioners can select the most appropriate approach for their specific requirements. Additionally, continuous exploration and integration of advanced techniques hold the promise of further enhancing the accuracy and robustness of language detection systems, paving the way for more effective cross-lingual communication and information processing.