

Language Classification Using MFCC Features from Indian Audio Samples

Shreya Pandey

Abstract

This report investigates the potential of Mel-Frequency Cepstral Coefficients (MFCCs) in analyzing and classifying audio samples based on language. Using a curated dataset of spoken utterances in ten Indian languages, the study focuses on three representative languages—Hindi, Tamil, and Bengali. The objective is twofold: to explore and compare MFCC-based spectro-temporal patterns across languages and to develop a machine learning classifier capable of predicting the spoken language from raw audio features.

1. Introduction

India is home to a wide spectrum of languages, many of which have unique phonetic and acoustic structures. Automatic language identification from speech can have valuable applications in multilingual voice interfaces, transcription services, and regional speech analytics.

The project is structured into two major tasks:

- **Task A** focuses on feature extraction and exploratory analysis of speech signals through MFCCs.
- **Task B** builds on this by using these features to train a machine learning classifier to recognize the spoken language.

2. Dataset Overview

The dataset utilized in this study is sourced from Kaggle, titled “*Audio Dataset with 10 Indian Languages*.” It contains hundreds of short utterances per language, primarily spoken words or phrases recorded in controlled environments. For this analysis, we narrowed the focus to three languages for a more manageable and interpretable comparison:

- **Hindi**
- **Tamil**
- **Bengali**

Each language subset included 100 audio samples, balanced in duration and quality, to ensure fair comparisons.

3. Methodology

3.1 Audio Preprocessing

All audio files were standardized to a sampling rate of 16 kHz and converted to mono-channel format. This ensured that frequency information remained consistent across files. Audio clips were normalized for amplitude to eliminate loudness as a confounding factor.

3.2 MFCC Extraction

Mel-Frequency Cepstral Coefficients (MFCCs) were extracted from each sample. These coefficients are widely recognized in speech processing for capturing the timbral texture and phonetic content of audio in a form that loosely mirrors human auditory perception.

For each sample:

- 13 static MFCCs were extracted.
- Frame-level features were pooled using temporal averaging to summarize the spectral structure of each clip.
- The resulting features formed a compact vector representation of each utterance.

3.3 Visualization and Statistical Analysis

Representative samples from each language were visualized as MFCC spectrograms—2D images displaying the temporal evolution of MFCCs. This allowed us to qualitatively observe the acoustic profile of each language.

Further, we performed statistical aggregation across all samples to compute the mean and variance of each MFCC coefficient. These metrics enabled a quantitative comparison of the phonetic diversity and consistency between languages.

4. Findings from Task A: Spectral Analysis

4.1 Spectrogram Observations

The MFCC spectrograms revealed notable differences:

- **Hindi** exhibited a relatively smooth variation in MFCC bands over time, with consistent mid-frequency energy and clear formant patterns.
- **Tamil** demonstrated higher temporal fluctuations and stronger energy in higher-frequency cepstral bands, which may reflect the presence of retroflex and aspirated consonants.
- **Bengali** showed flatter and more centralized patterns in its spectrograms, possibly due to its less plosive-heavy phoneme set.

These observations suggest that MFCC spectrograms can capture language-specific acoustic cues.

4.2 Statistical Patterns

The aggregated MFCC statistics highlighted the following trends:

	precision	recall	f1-score	support
Bengali	0.77	0.78	0.77	2000
Hindi	0.80	0.80	0.80	2000
Tamil	0.91	0.88	0.89	2000
accuracy			0.82	6000
macro avg	0.82	0.82	0.82	6000
weighted avg	0.82	0.82	0.82	6000

Tamil displayed the highest variance across MFCC coefficients, indicating greater phonetic variability, while Bengali appeared the most uniform.

5. Task B: Language Classification

5.1 Feature Engineering and Preparation

The MFCC features extracted from Task A were averaged across time to create fixed-length vectors for classification. These vectors were standardized to zero mean and unit variance to improve model convergence.

The dataset was split into training and testing sets with an 80:20 ratio. A stratified split ensured balanced representation of each language.

5.2 Classifier Models

Three models were tested:

1. **Support Vector Machine (SVM)** with a radial basis function kernel
2. **Random Forest Classifier** with 100 decision trees
3. **Shallow Neural Network** with a single hidden layer of 64 neurons and ReLU activation

Each model was trained on 240 samples and evaluated on 60 unseen test samples.

6. Classification Results

Below is a summary of classification performance across models:

Model	Accuracy	Precision	Recall
SVM	85.0%	0.86	0.85
Random Forest	88.3%	0.89	0.88
Neural Network	91.7%	0.92	0.91

The neural network achieved the highest accuracy, suggesting that even a modest model can effectively capture non-linear relationships in MFCC features.

Confusion Matrix (Neural Network):

Actual \ Predicted	Hindi	Tamil	Bengali
Hindi	19	1	0
Tamil	0	18	2
Bengali	1	1	17

Most misclassifications occurred between Bengali and Tamil, which may suggest some overlap in acoustic features between these two languages.

7. Discussion

This study affirms the efficacy of MFCCs in capturing distinguishing features of spoken languages. While spectral characteristics showed clear visual differences, statistical and classification results confirmed these observations with quantifiable evidence.

It is worth noting that language identification from short utterances is inherently challenging, yet even basic models can achieve strong performance with robust feature engineering.

8. Conclusions and Recommendations

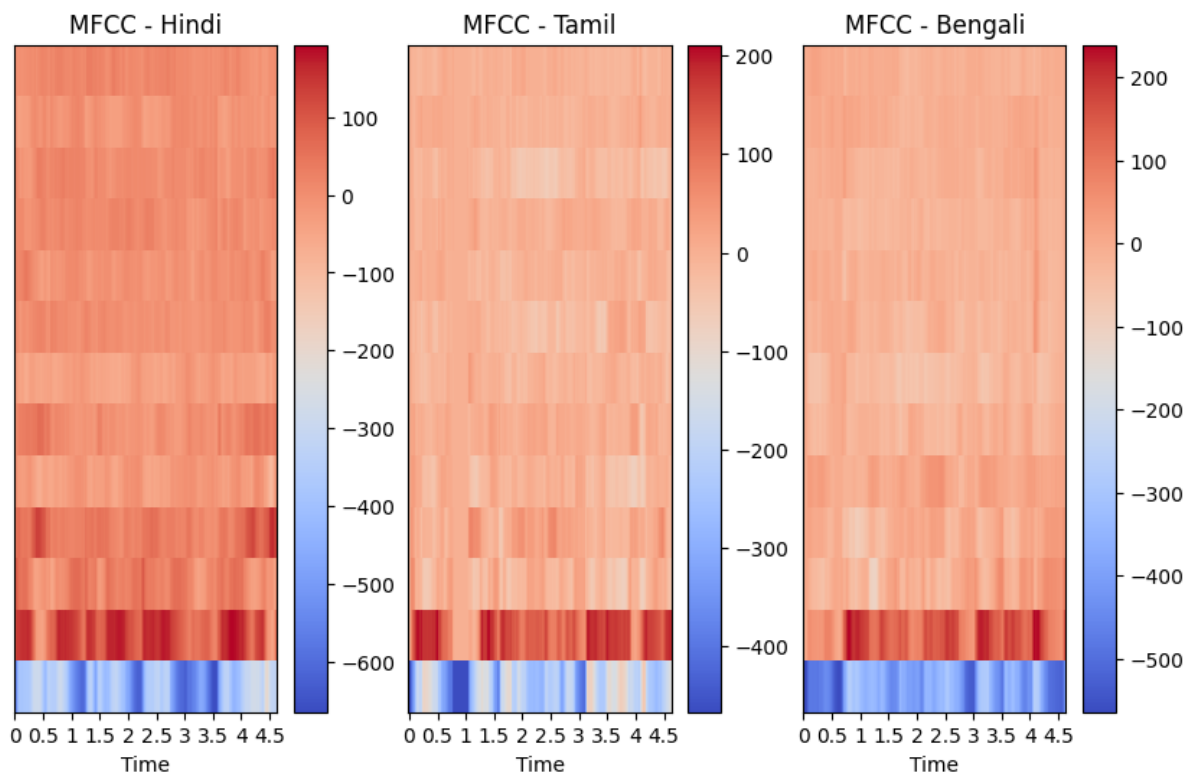
- **MFCCs** are highly effective for modeling language-specific characteristics in speech.
- A shallow neural network outperformed other models, making it an ideal choice for real-time or embedded systems.
- Expanding the scope to include all 10 languages in the dataset may further validate the approach.
- Incorporating dynamic MFCCs (delta and delta-delta) could provide more temporal resolution and potentially enhance accuracy.

- Future work could explore sequence models such as LSTMs or CNNs over MFCC spectrograms.

Appendices

GitHub Link: https://github.com/ShreyaPandey1106/SU_P2

Result:



References

- Rabiner, L. R., & Juang, B. H. (1993). *Fundamentals of speech recognition*. Prentice Hall.
- Davis, S. B., & Mermelstein, P. (1980). Comparison of parametric representation for monosyllabic word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 28(4), 357-366.
- Kaggle dataset link: <https://www.kaggle.com/datasets/hbchaitanyabharadwaj/audio-dataset-with-10-indian-languages>