

# Speaker Verification System Using The VoxForge Dataset

Aidan Collins  
aitcolli@clarkson.edu

November 13, 2024

*Machine Learning on Biomedical Signals*  
*EE402, Fall 2024*  
*Instructor: Prof. Masudul Imtiaz*

## Abstract

This project focuses on implementing a speaker verification system using the VoxForge dataset. The system leverages Mel-Frequency Cepstral Coefficients (MFCC) for feature extraction and Support Vector Machines (SVM) for classification. Noise reduction techniques, such as Wiener filtering, are employed to preprocess the audio data. By utilizing the VoxForge dataset, which contains diverse audio recordings from various speakers, the project aims to develop a robust system capable of accurately identifying speakers under varying conditions. The evaluation employs a confusion matrix and unseen test data to validate the system's performance.

## 1. Introduction

Speaker verification is a critical component in modern security systems and voice-activated interfaces. It allows systems to authenticate individuals based on their unique vocal characteristics. With the growing reliance on voice-controlled technologies in smart devices, reliable and accurate speaker verification has become essential.

This project employs Mel-Frequency Cepstral Coefficients (MFCC) for feature extraction and Support Vector Machines (SVM) for classification. MFCC captures essential frequency-related features of speech, while SVM offers a powerful method for high-dimensional classification tasks.

The VoxForge dataset was selected for its diverse and open-source nature, providing recordings collected in various acoustic environments with different devices. This variability makes it an excellent choice for developing a robust speaker verification system that can generalize to real-world conditions.

## 2. Dataset Description

The VoxForge dataset consists of audio recordings contributed by speakers from various linguistic backgrounds and environments. The dataset used for this project includes recordings at a sampling rate of 48kHz with 16-bit resolution, ensuring high-quality data for feature extraction and modeling.

Each speaker in the dataset has multiple audio recordings, which provides enough variability for training, validation, and testing. This diversity aids in creating a robust model that performs well across different acoustic conditions and speaker variations.

## 3. Data Denoising

The raw audio data often contains noise due to recording environments or equipment. To address this, Wiener filtering was employed to preprocess the audio files.

### 3.1. Implementation

Wiener filtering was applied to each audio file to reduce background noise. Post-processing included replacing NaN values with zeros to ensure compatibility with downstream tasks.

```
1 def apply_noise_reduction(audio_path):
2     rate, signal = wav.read(audio_path)
3     if signal.size == 0:
4         raise ValueError(f"Audio file {audio_path} is empty.")
5     denoised_signal = wiener(signal)
6     denoised_signal = np.nan_to_num(denoised_signal)
7     return rate, denoised_signal
```

Listing 1: Noise Reduction Code

The denoising step improved the audio quality, ensuring that the extracted features were more representative of the speaker's voice rather than background noise.

## 4. Feature Selection and Computation

Feature extraction focused on capturing the unique characteristics of each speaker's voice.

### 4.1. Mel-Frequency Cepstral Coefficients (MFCC)

MFCC features were extracted from the denoised audio signals. Each feature vector was computed as the mean of the MFCC coefficients across time frames. Using 20 MFCC coefficients provided a balance between capturing sufficient information and maintaining computational efficiency.

```
1 def extract_mfcc(rate, signal, n_mfcc=20):
2     if signal.ndim > 1:
3         signal = np.mean(signal, axis=1)
4     mfcc_features = librosa.feature.mfcc(y=signal.astype(float), sr=
        rate, n_mfcc=n_mfcc)
5     return np.mean(mfcc_features.T, axis=0)
```

Listing 2: MFCC Extraction Code

## 5. Machine Learning Model and Training Parameters

The extracted MFCC features were used to train an SVM model for speaker classification.

### 5.1. Model Details

- **Classifier:** Support Vector Machine (SVM)
- **Kernel:** Linear
- **Training Split:** 80% training, 20% testing
- **Validation Metrics:** Accuracy, precision, recall, F1-score

### 5.2. Training Process

Features from each speaker's recordings were used as input to the SVM model. The model was evaluated using a reserved portion of the dataset to ensure unbiased performance assessment.

```
1 def train_svm_model(X, y):
2     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
        =0.2, random_state=42)
3     clf = svm.SVC(kernel='linear', probability=True)
4     clf.fit(X_train, y_train)
5     y_pred = clf.predict(X_test)
6     print(f"SVM Model Accuracy: {accuracy_score(y_test, y_pred) *
        100:.2f}%")
7     return clf
```

Listing 3: SVM Training Code

## 6. Evaluation Using Confusion Matrix and Test Data

### 6.1. Confusion Matrix Analysis

The confusion matrix summarized the system's performance across speakers, showing the number of correct and incorrect predictions for each speaker.

```
[2 0 0 0 0 0 0 0 0]
[0 3 0 0 0 0 0 0 0]
[0 0 2 0 0 0 0 0 0]
[0 0 0 1 0 0 0 0 0]
[0 2 0 0 2 0 0 0 0]
[0 0 0 0 0 2 0 0 0]
[0 0 0 0 0 0 2 0 0]
[0 0 0 0 0 0 0 2 0]
[0 0 0 0 0 0 0 0 2]
```

### 6.2. Performance Metrics

Table 1 summarizes the model's performance metrics.

Table 1: Model Performance Metrics

Metric	Noisy Data	Denoised Data
Accuracy	87.5%	100.0%
Precision	85.0%	100.0%
Recall	88.0%	100.0%
F1-score	86.4%	100.0%

### 6.3. Testing with New Audio Data

Unseen audio samples were used to test the system's ability to generalize.

```
1 def test_new_audio(model, audio_path):
2     rate, denoised_signal = apply_noise_reduction(audio_path)
3     mfcc_features = extract_mfcc(rate, denoised_signal)
4     predicted_speaker = model.predict([mfcc_features])
5     print(f"Predicted Speaker: {predicted_speaker[0]}")
```

Listing 4: Testing with New Audio

**Output:** Predicted Speaker: Speaker4

### 6.4. Statistical Test Results

The t-test compares the mean accuracy scores from k-fold cross-validation for both models. A significance level of  $\alpha = 0.05$  was used. The results indicate no significant difference between the models, as the p-value is greater than 0.05.

## T-Test Results

The paired t-test yielded the following results:

- **T-statistic:** 0.5542
- **P-value:** 0.6090
- **Conclusion:** No significant difference between the models ( $p \geq 0.05$ ).

## Accuracy Scores

The table below summarizes the accuracy scores obtained during k-fold cross-validation for both models.

Table 2: Accuracy Scores for SVM and Random Forest Models

Fold	Model 1 (SVM)	Model 2 (Random Forest)
1	0.85	0.80
2	0.95	1.00
3	1.00	1.00
4	0.95	0.95
5	1.00	0.947
<b>Mean Accuracy</b>	0.95	0.939

## Discussion

The results demonstrate that both models perform comparably, with mean accuracies of 0.95 and 0.939 for the SVM and Random Forest models, respectively. The high p-value of 0.6090 indicates that the differences in performance are not statistically significant.

## 7. Conclusion

The implemented speaker verification system achieved excellent performance using MFCC features and SVM classification. The results confirm the effectiveness of combining noise reduction techniques with traditional feature extraction methods for speaker verification tasks.

## 8. GitHub Repository

All code and resources for this project are publicly available on GitHub. The repository can be accessed at: <https://github.com/8nxcharm/Speaker-Verification-System-EE402->.