

Speech Graphics - Speech Mode Classification

Sambit Paul

September 24, 2018

Abstract

This report outlines the project done for speech mode classification. The idea is to use deep learning methods to classify audio files for types of speech modes - rap, singing and speech. The deep learning models used here are feed-forward networks, LSTM-Recurrent Networks and Convolutional Neural Networks. Additionally, traditional methods like K-nearest neighbours and SVMs are also used primarily to compare the performance of deep learning methods against traditional methods. The features used for almost all the models were the MFCC features except for the CNN which used the spectrogram of the audio files. From the experiments, it was observed that the LSTM-RNN using the MFCC features was able to classify with the highest accuracy.

1. Problem Description

The objective is to classify audio segments on basis of three modes — rap, singing and speech. The dataset provided here contains a total of 1440 audio recording each of which has a length of 3 seconds and are recorded at 16kHz. Each mode of speech has 480 examples associated to them.

The idea is to extract features from the audio files and use them for classifying between the three given modes. The data needs to be distributed such that 10% of the dataset is in a balanced test set.

2. Data

This section explores the data to be used for the experiments. The first section explores features that can be used for the experiments. The second section aims to visualise the features of the data from the high dimensional space projected to a 2 dimensional space.

2.1. Feature Extraction

One of the key challenges for speech recognition is selecting the optimal feature set which will generate the best results. The problem, as stated in Section 1, is a basic problem in speech classification which will work well given the same features as most problems in Automated Speech Recognition.

The first part discusses about Mel-Frequency Cepstral Coefficients (MFCC) and Linear Prediction Cepstral Coefficient (LPCC) and in the second part we look into spectrograms and their potential use in this classification problem.

MFCC is based on human hearing such that the filters below 1kHz are linearly spaced and the ones above 1kHz are logarithmically spaced. MFCC operates by converting the signal from time domain to frequency domain using Fourier Transform and then, each Mel filter computes a weighted sum of its spectral components. Discrete Cosine Transform (DCT) applied to the log Mel spectrum produces the MFC coefficients [Muda et al., 2010].

LPCC works on the principle that a speech sample at the current time can be determined by a linear combination of previous speech samples. LPCC operates by applying Linear Predictive Analysis on the speech signal followed by a Cepstrum Analysis which finds the cepstrum coefficients of the speech sequence [Gulzar et al., 2014].

MFCC is considered better than LPCC for most problems in Automatic Speech Recognition [Rana and Miglani, 2014]. Figure 1 shows MFCC samples taken from each speech mode with 13 cepstrums.

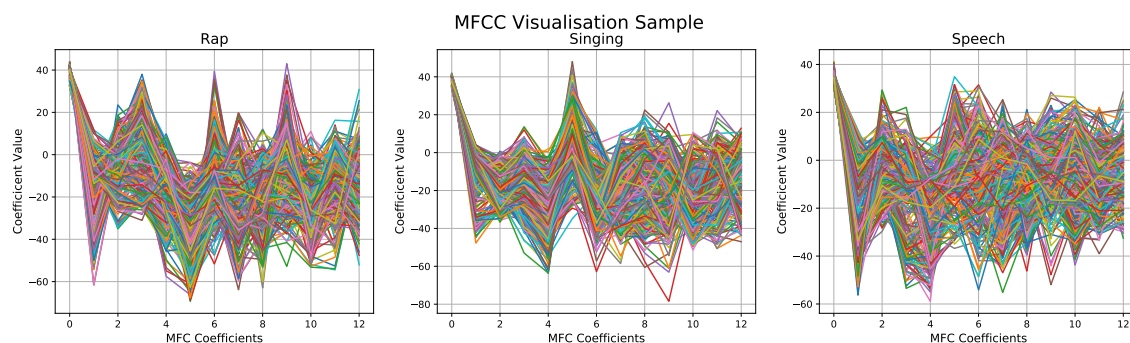


Figure 1: Visualisation of MFCC for each speech mode. Each plot contains 299 lines each representing a frame.

One of the features that could potentially also be useful is the spectrogram. There is a significant difference in the spectrograms of the different speech modes. Figure 2 shows sample spectrograms of the different speech modes. It can be the case, the singing sample has a dominant higher frequency throughout the sample compared to other modes, while speech has gaps at some intervals in the frequency map.

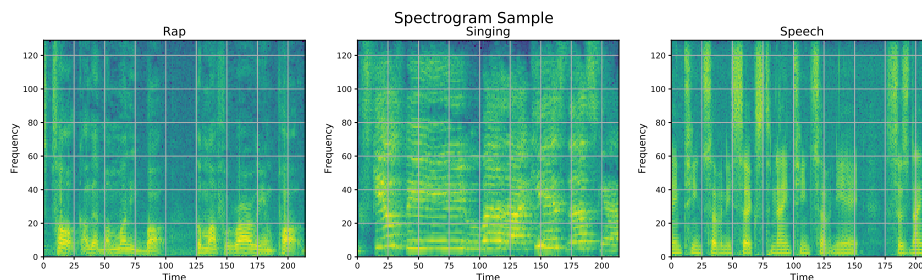


Figure 2: Sample spectrogram for each speech mode

2.2. Visualisation with t-SNE

This section aims to visualise the MFCC features using T-distributed Stochastic Neighbour Embedding (t-SNE) to see if any clusters are formed when projecting the data into two dimensions.

t-SNE is a tool for visualising high dimensional data [Maaten and Hinton, 2008]. We are trying to project each audio file from their high-dimensional MFCC values to 2 dimensional space. Due to the high-dimensionality of the data ($299 \text{ frames} \times 13 \text{ cepstrums}$ for each file), we are using Principal Component Analysis (PCA) to reduce dimensions to limit computations. We use 200 components in the PCA which has an explainability factor of 77.13%.

Figure 3 shows projections of the data for different perplexities but as is visible, there are no distinct clusters visible for any value of perplexity.

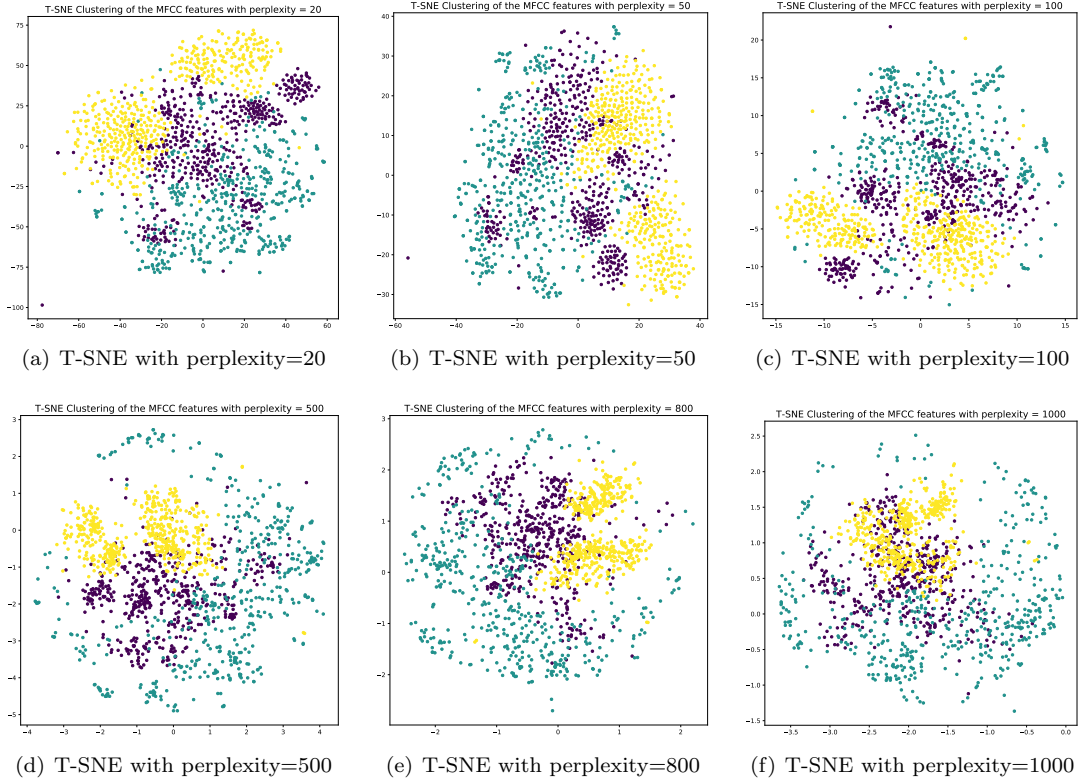


Figure 3: T-SNE projections of the dataset for perplexity values of 20, 50, 100, 500, 800 and 1000

3. Deep Learning Models

This section explores different deep learning models and their performance compared to traditional machine learning models.

For these experiments, we break down the dataset into three parts — the training set which contains 1101 examples (76.5%), the validation set which contains 195 examples (13.5%) and the test set which contains 144 examples (10%).

For the feed-forward network and LSTM-RNN, we use the MFCC features and for the CNN, we use images of the spectrogram as inputs. The loss function for all the models is the Cross Entropy loss which means that the class with the highest energy value is the most probable class.

3.1. Feed-Forward Network

This section describes the feed-forward model used for classification using a flattened form of the 299×13 dimensional MFCC. The architecture has 5 fully connected layers starting with 3887 (input layer) followed by 1024, 256, 64, 8 and finally 3 output nodes each holding the cross-entropy energies of each class. Figure 4 shows the feed-forward architecture used here.

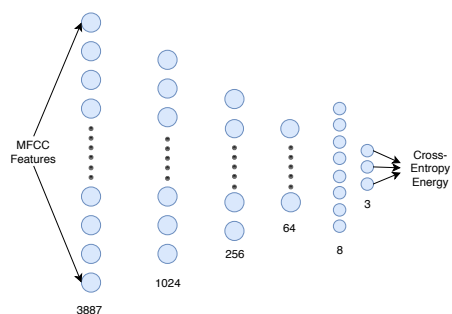


Figure 4: Feed-Forward Network Architecture

The training-validation loss curve for the feed-forward architecture is shown in Figure 5.

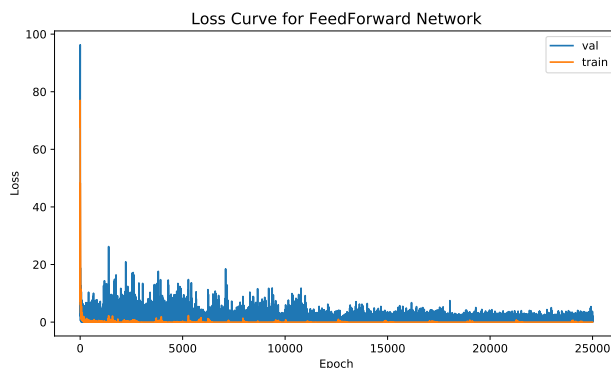


Figure 5: Training-Validation Loss Curve

The test set accuracy for this model was 84.03%.

3.2. Convolutional Neural Network

This section describes the Convolutional Neural Network (CNN) used for classification using the spectrogram images. On investigation of the spectrograms of different speech modes, it can be observed that each category has a specific nature of associated to it which the CNN should be able to extract as features. The model has 5 convolutional layers followed by a fully-connected 5 dense layers which gives the cross entropy energies of each class at the end. The kernel size of the convolutional layers is 5×5 in the first 2 layers followed by 3×3 in consequent convolutional layers maintaining a constant stride of 1. A detailed architecture of the network is given in Figure 6.

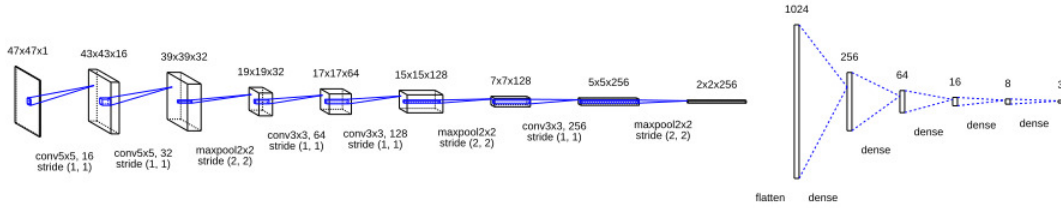


Figure 6: CNN Architecture

The training-validation loss curve for the feed-forward architecture is shown in Figure 7. It can be observed, compared to the feed-forward model loss (Figure 5), the loss values have a much smaller variance. This can be used to infer that this model should have a better performance.

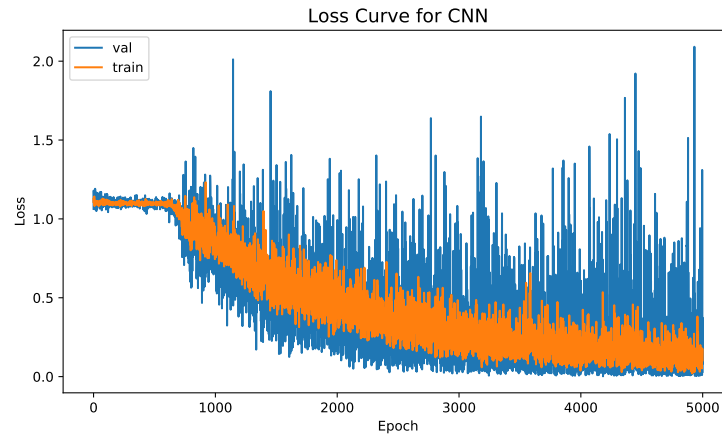


Figure 7: Training-Validation Loss Curve

The test set accuracy for this model was 88.19%.

3.3. Long Short-Term Memory RNN

This section describes the LSTM-Recurrent Neural Network used for classification using the MFCC features. We have used 2 LSTM layers in this model followed by a 2 dense layers. The LSTM takes 13 dimensional input with a sequence length of 299. The architecture of the 2 layer LSTM network is given in Figure 8.

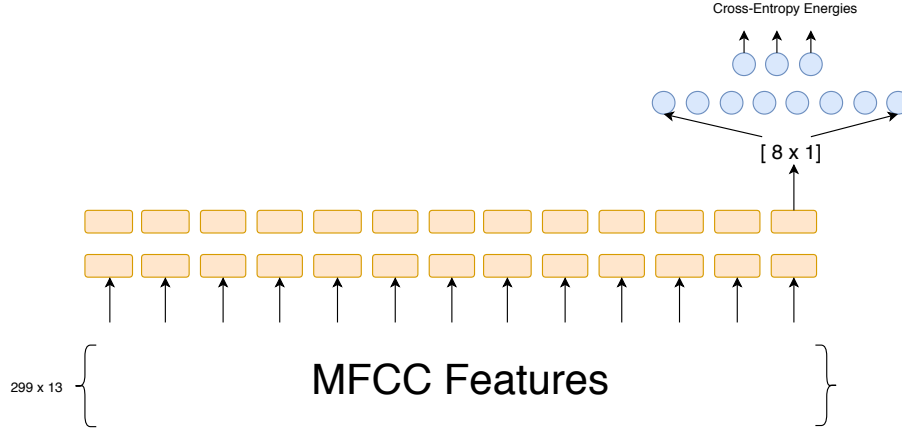


Figure 8: LSTM Architecture

The training-validation loss curve for the LSTM-Recurrent Network is shown in Figure 9. Comparing to the CNN performance as shown in Figure 7, it can be observed that the losses are almost similar and hence, the performance can also be expected to be nearly same.

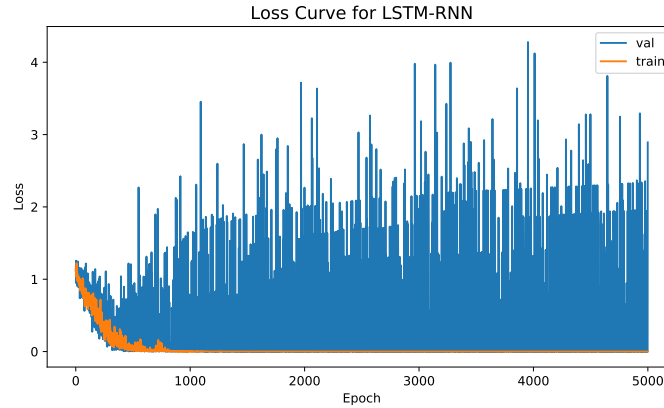


Figure 9: Training-Validation Loss Curve

The test set accuracy for the 2-layer LSTM model was 89.58%.

4. Traditional Methods

In addition to the deep learning methods, experiments were also performed using traditional machine learning methods to carry out a comparison between traditional and deep models. For these experiments, we break the dataset into 1296 examples (90%) in the training set and 144 examples (10%) in the test set. The features used are the MFCC features as extracted using the python-speech-features library [Lyons, 13].

4.1. K-Nearest Neighbours

The first traditional method used is the uniformly weighted K-nearest neighbours classifier. The experiment used 5 neighbours which yielded a test accuracy of 63.19%. The low performance of k-NN may be alluded to the fact that, since the dimensions of the data is so high, the notion of distance between different data points is quite blurred.

4.2. Support Vector Machine

The second method used support vector machines based on Linear Support Vector Classification with L2 regularisation. We do not use L1 norm since the sparse nature of the resulting coefficient matrix might lose out on a few of the dimensions while generating the hyperplane which we do not want. This method yielded a test accuracy of 71.53%.

5. Results & Analysis

Table 1 gives a summary of the test performance of all the methods used.

Method	Test Accuracy
Feed Forward Network	84.03%
Convolutional Neural Network	88.19%
LSTM-RNN-1	86.11%
LSTM-RNN-2	89.58%
KNN-5	63.19%
SVM-LinearSVC	71.53%

Table 1: Test accuracies of the different methods explored in this project.

As is visible, traditional methods like KNN and Linear SVC do not perform at par with deep learning methods. This implies that for this dataset, the classification boundaries are not linear and therefore, introduction of the non-linear activations in the deep networks helped achieve much better performance. It might be possible to achieve better performance with SVMs if non-linear kernels like Radial Basis Functions or Polynomial Functions were used [Aggarwal and Aggarwal, 2011].

Moreover, the fact that LSTM networks generated $\sim 5\%$ improvement over the feed-forward network reveals the importance of exploiting the sequential nature of the data.

Although, we can observe that LSTMs have performed better than CNNs by $\sim 1\%$, further analysis of hyper-parameters and architecture needs to be done to see if CNNs with spectrogram inputs can outperform LSTM networks.

6. Real-Time Classification

For real-time classification, using an LSTM model would be of useful since they do not require fixed length inputs. If audio is recorded at 16kHz, we would get 16000 datapoints which when passed through MFCC algorithm (frame size = 0.025, frame step = 0.01, number of cepstrum = 13), should return a 98×13 feature matrix based on Equation 1.

$$n_{frame} = \frac{|seq_{len} - (frame_{size} * SR)|}{(frame_{stride} * SR)} \quad (1)$$

where,

$n_{frame} = \text{Numberof frames}$

$seq_{len} = \text{SequenceLength}$

$frame_{size} = \text{FrameSize}$

$frame_{stride} = \text{FrameStride}$

$SR = \text{SamplingRate}$

Another technique that would be interesting to try out would be a recurrent convolutional network which can take in spectrograms generated in real-time and find out correlations between subsequent spectrograms. It will potentially be able to learn how transitions between different speech modes work, which can be critical for classification.

References

- [Aggarwal and Aggarwal, 2011] Aggarwal, S. and Aggarwal, N. (2011). Classification of audio data using support vector machine 1. *International Journal of Computer Science and Technology*.
- [Gulzar et al., 2014] Gulzar, T., Singh, A., and Sharma, S. (2014). Comparative analysis of lpcc, mfcc and bfcc for the recognition of hindi words using artificial neural networks. *International Journal of Computer Applications*, 101(12):22–27.
- [Lyons, 13] Lyons, J. (2013–). python-speech-features: Provides common speech features for asr including mfccs and filterbank energies for Python.
- [Maaten and Hinton, 2008] Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.
- [Muda et al., 2010] Muda, L., Begam, M., and Elamvazuthi, I. (2010). Voice recognition algorithms using mel frequency cepstral coefficient (mfcc) and dynamic time warping (dtw) techniques. *arXiv preprint arXiv:1003.4083*.
- [Rana and Miglani, 2014] Rana, M. and Miglani, S. (2014). Performance analysis of mfcc and lpcc techniques in automatic speech recognition. *International Journal Of Engineering And Computer Science*, 3(08).