

# Speaker Accent Recognition by MFCC Using K-Nearest Neighbour Algorithm: A Different Approach

Munish Bhatia<sup>1</sup>, Navpreet Singh<sup>2</sup>, Amitpal Singh<sup>3</sup>

Student, Computer Science, GNDU, Gurdaspur, India<sup>1,2</sup>

Assistant Professor, Computer Science, GNDU, Gurdaspur, India<sup>3</sup>

**Abstract:** A K-Nearest Neighbour Algorithm involving Mel-Frequency Cepstral Coefficients (MFCCs) is provided to perform Speech signal feature extraction for the task of speaker accent recognition. Mel-Frequency Cepstral Coefficient is effectively used to perform the feature extraction of the input signal. For each input signal the mean of the MFCC matrix is used for pattern recognition. The K-nearest neighbour algorithm is based on evaluating minimum Euclidean distance measure from input data set to stored data set. Since large number of speakers of different accent are present, they can be grouped together depending upon their accent. Thus each signal coming from different group makes a distinct MFCC vector. In this paper we have compared the MFCC from global group to smaller sub groups.

**Keywords:** MFCC, K-nearest neighbour, Voice Code, Euclidean distance, Global group

## I. INTRODUCTION

A speaker's accent distinguishes him/her as a member of a group. These groups have been classified on the basis of geographic areas, by social class or by ethnicity. As human being we are sensitive to accent and can tell whether a speaker belongs to a group or not. For above scenario to be possible, accent are pattern of speaking that distinguish a group from other group. It is assumed that these stylised patterns affect word rate, the phonologic coding used in lexicon and the prosody of utterances [8][10]. Speech technology is still facing problem with pronunciation variation across different accent groups. In automated speech recognition, a variation in accent used in testing and training can lead increase in word error rate. In synthesis of speech, voices are constant in one accent due to large use of corpus based methods operating by single speaker speech [9].

Speech Processing involves capturing of sound and process its features by extracting it. The technique for extracting feature is done by MFCC i.e. Mel Frequency Cepstral Coefficient [1]. MFCC is by far the most successful technique for performing feature extraction. MFCC is depiction of short term power spectrum on the basis of linear cosine transform of log spectrum on non-linear Mel scale of frequency.

## II. STEPS FOR COMPUTING MFCC

### 1. Pre-emphasis Filtering

The input speech signal  $q(n)$  is sent to high pass filter

$$q_2(n) = q(n) - a \cdot q(n-1)$$

Where  $q_2(n)$  is the output speech signal and  $a$  is valued between 0.9 and 1.0. The z-transform is given by

$$T(z) = 1 - a/z$$

The goal is to compensate the high frequency portion that was suppressed during production of sound.

### 2. Framing

The speech signal then is segmented into frames of ~25 ms with overlap of ~1/3 frame size.

### 3. Hamming Window

Each frame is multiplied with hamming window for the continuity of the first and last point of frame. Hamming windowing is  $q(n) \cdot w(n)$  where  $w(n)$  is the hamming window.

### 4. Fast Fourier Transform

In this step we perform the FFT to get the magnitude of the frequency of each frame assuming that signal inside frame is periodic.

### 5. Triangular Band pass Filtering

Then magnitude of the frequency is multiplied by set consisting of about triangular band pass filter. The positioning of the filter is synchronous with the Mel frequency. Mel frequency is proportional to log of linear frequency.

### 6. Discrete Cosine Transformation

In this step DCT is applied to N log Energy resulting from band pass filter. Usually  $N=20$ . The result is the features similar to cepstrum that's why these are termed as Mel scale Cepstral Coefficients or MFCC [5].

Along with above, other dynamic features like Del1 (First order derivative) and Del2 (Second Order derivative) are included. As of now 13 MFCC coefficients are computed along with Del1 and DEL2 making total up to 39 coefficients. One important point to be noted that for each signal the MFCC actually form  $n \cdot p$  matrix as shown in figure 1 where  $n$  is the number of window frames and  $p$  is the number of MFCC coefficient i.e. 39. Now the MFCC is passed to pattern recognition techniques [2][3] and for that these matrices are converted to vectors by simply taking the mean of the  $n$  column vectors. Fig 1 shows an instance of MFCC matrix. Usually as mentioned  $N=20$ . The mean of the column vector is shown in Fig 2. Mean vector is calculated as

$$a_1 = (x_1 + y_1 + \dots + z_1) / n$$

$$a_2 = (x_2 + y_2 + \dots + z_2) / n \text{ and so on.}$$

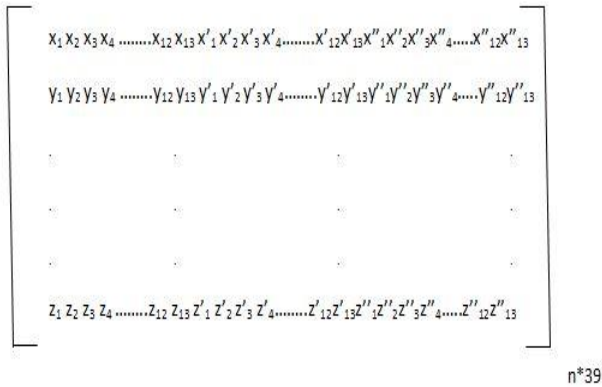


Fig1.MFCC matrix

Mean vector is considered as Voice code (V\_C). The three partition of the vector are termed as VOICE CODE 1 (V\_C1), VOICE CODE 2 (V\_C2), and VOICE CODE 3 (V\_C3) respectively.

$$A = [a_1 a_2 a_3 a_4 \dots a_{12} a_{13} b_1 b_2 b_3 b_4 \dots b_{12} b_{13} c_1 c_2 c_3 c_4 \dots c_{12} c_{13}]$$

Fig2 Mean Voice Code MFCC Vector

In Fig2. A is the Mean Voice Code Vector, V\_C1 is from a<sub>1</sub> to a<sub>13</sub>, V\_C2 is from b<sub>1</sub> to b<sub>13</sub>, V\_C3 is from c<sub>1</sub> to c<sub>13</sub>.

**K -NEAREST NEIGHBOUR ALGORITHM FOR ACCENT RECOGNITION**

*Basic Assumption*

In order to assign a member to a group, initially it is assumed that all the different accent members form a single group which is considered as global group. Within that global group different accent are grouped to form local group with closely similar accents [4]. Now these sub groups are further partitioned into sub subgroups called regions which are different from each other. The number of partitions can be varied. Above scenario can be shown in Fig 3

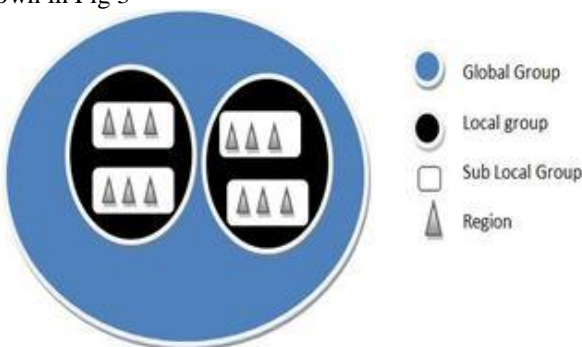


Fig3. Grouping of closely related Members in Sub Groups

*Step 1*  
Initially Input signal's V\_C1 compared for Euclidian distance as given by Equation 1, with each mean of V\_C1 of different local groups of Fig 4 showing global Group C with different sub Groups like C1, C2, C3 etc. These sub groups contain the mean of V\_C1 of the stored set.

Equation 1  
N 13

$$L\_DIST = \min_{i=1} \sum_{k=1} (|A_k - Z_{ik}|)$$

Where L\_DIST is the minimum local distance of the global group to local group, A<sub>k</sub> is the input signal V\_C1 and Z is the stored data set and N is the number of stored data set to compare with. Now as soon as the minimum distance is calculated, it can be estimated that the input signal is closely related the particular local group of minimum global distance. For instance as shown in Fig 4.

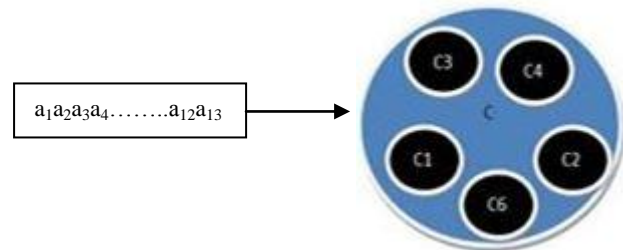


Fig4. Input Signal V\_C1 compared with mean V\_C1 of local groups in C

Let us assume that the input signal V\_C1 is closely related to mean V\_C1 of local group C1.

*Step 2*

After determining the local group, the next step is to determine the sub local group to which the member belongs. For that the comparison of input V\_C2 is done with the sub local groups containing their respective V\_C2's mean as shown in Fig5. Comparison using Equation 2 gives the minimum distance of local group to sub local group. Variables have their usual meanings.

Equation 2  
N 26

$$SL\_DIST = \min_{i=1} \sum_{k=14} (|A_k - Z_{ik}|)$$

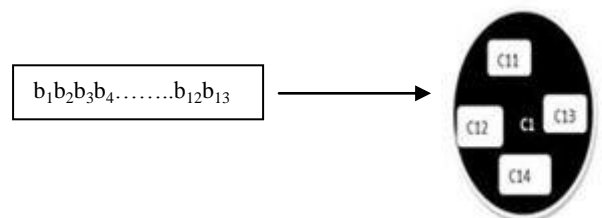


Fig5. Input signal V\_C2 compared with mean V\_C2 of Sub local Groups in C1

Assume that the input V\_C2 is closely related to V\_C2 of Sub local group C11.

*Step 3*

Similarly as above, repeating the step one more time we get the final region to which the member belongs. Finally the V\_C3 is compared according to equation 3 which

determines the minimum region distance and the final region of the member with the minimum distance given by  $DIST=L\_DIST+SL\_DIST+R\_DIST$

Equation3  
N 39

$$R\_DIST = \min_{i=1} \sum_{k=27} (|A_k - Z_{ik}|)$$

Fig 6 shows the final step.

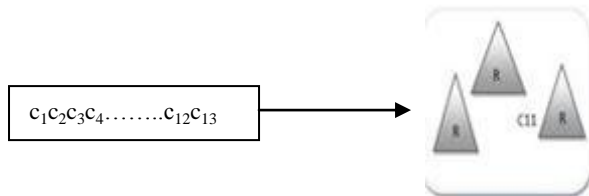


Fig6. Input signal V\_C3 compared with mean V\_C3 of regions in C11

### CONCLUSION

We report the different approach of accent determination by using the K-nearest neighbour algorithm based on MFCC. In this paper we reviewed the classification of accents on the basis of MFCC coefficient obtained from input signal. We also considered the training data must be kept for determination of accent to greater accuracy. Further we are planning to compare different classifiers for different accents.

### ACKNOWLEDGMENT

We would like to thank GNDU facility for all the help they have provided.

### REFERENCES

- [1] Liu Wai Kat Fung, P. "Acoustics, Speech, and Signal Processing, 1999", Proceedings 1999 IEEE International Conference.
- [2] Zichen Ma, Ernest Fokoué "A Comparison of Classifiers in Performing Speaker Accent Recognition Using MFCCs" published Online June 2014 in SciRes.
- [3] Sandesh Aryal and Ricardo Gutierrez-Osuna "Accent Conversion through Cross Speaker Articulation" in IEEE proceedings, ICASSP 2014
- [4] Malay K. Pakhira "A Modified k-means Algorithm to Avoid Empty Clusters" in International Journal of Recent Trends in Engineering, Vol 1, No. 1, May 2009
- [5] MFCC Coefficients Available at <http://www.mirlab.org/jang/books/.../speechFeatureMfcc.asp?title=12-2%20MFCC>.
- [6] Gina-Anne Levow, "Investigating Pitch Accent Recognition in Non-native Speech".
- [7] Carlos Teixeira, Isabel Trancoso and Antonio Serralheiro. "Accent Identification".
- [8] M Huckvale, "ACCDIST: a metric for comparing speakers' accents", International Conference on Spoken Language Processing 2004.
- [9] Chao Huang, Eric Chang, and Tao Chen "Accent Issues in Large Vocabulary Continuous Speech Recognition (LVCSR)", Microsoft Research.
- [10] Fadi Biadsy, "Automatic Dialect and Accent Recognition and its Application to Speech Recognition".