

Improving the Noise-Robustness of Mel-Frequency Cepstral Coefficients for Speech Processing

Sourabh Ravindran¹, David V. Anderson¹, Malcolm Slaney²

¹School of Electrical and Computer Engineering
Georgia Institute of Technology, Atlanta GA 30332

²Yahoo! Research

701 First Avenue, Sunnyvale, CA 94089

sourabh@gatech.edu, dva@ece.gatech.edu, malcolm@ieee.org

Abstract

In this paper we study the noise-robustness of mel-frequency cepstral coefficients (MFCCs) and explore ways to improve their performance in noisy conditions. Improvements based on a more accurate model of the early auditory system are suggested to make the MFCC features more robust to noise while preserving their class discrimination ability. Speech versus non-speech classification and speech recognition are chosen to evaluate the performance gains afforded by the modifications.

1. Introduction

MFCC's are very useful features for audio processing in clean conditions. However, performance using MFCC features deteriorates in the presence of noise. There has been an increased effort in recent times to find new features that are more noise robust compared to MFCCs. Features such as, spectro-temporal modulation features [1] are more robust to noise but are computationally expensive. Skowronski and Harris [2] suggested modification of MFCC that uses the known relationship between center frequency and critical bandwidth. They also studied the effects of wider filter bandwidth on noise robustness. Herein, we suggest different modifications to MFCCs that make it more robust to noise without adding prohibitive computational costs.

MFCC features approximate the frequency decomposition along the basilar membrane by a short-time Fourier Transform. The auditory critical bands are modeled using triangular filters, compression is expressed as a log function and a discrete cosine transform (DCT) is used to decorrelate the features [3].

In this paper we cite two reasons for the poor noise performance of MFCCs. First, block processing with Fourier transform and the use of triangular filters for grouping of frequency bins into critical bands results in a signal in each channel that is not as smooth as that obtained with band-pass filters. We introduce a spatial subtraction stage to improve the performance of band-pass filter based features. Second, MFCC's poor noise performance can be attributed to log compression [4], [5], [6]. The large negative excursions of the log function for values close to zero leads to a splattering of energy after the DCT whereas root compression (expressed as $(\cdot)^\alpha$, with $0 < \alpha < 1$) followed by the DCT leads to better compaction of energy, as is shown later.

The rest of the paper is organized as follows, Section 2 explains the experimental setup. Sections 3, 4 and 5 talk about meth-

ods to improve noise-robustness of MFCCs. Section 6 uses an information theoretic measure of clustering to support the results obtained in earlier sections, followed by conclusions.

2. Experimental Setup

In order to study the noise-robustness of MFCCs, we choose a speech versus non-speech classification task at various signal-to-noise ratios (SNR). The results are also validated by performing speech recognition on the Aurora 2 database. The audio database for speech versus non-speech classification task was built from five publicly available corpora. Speech samples were taken from TIMIT acoustic-phonetic continuous speech corpus. The training set consisted of 300 examples from TIMIT's training subset and test set consisted of 150 different sentences spoken by 50 different speakers (25 male, 25 female) from TIMIT's test subset. Sentences and speakers in the training and test sets are different. The non-speech class consisted of 450 examples (300 for training and 150 for testing) which included animal vocalization from BBC Sound Effects audio CD collection, music samples from RWC Genre Database [7] and environmental sounds from NoiseX and Aurora databases. Segments of one second duration were selected for training and testing. The task consisted of predicting whether a given one second test segment belongs to the speech or the non-speech class. Pink noise was synthetically added to generate different SNRs.

Each audio segment was divided into frames of length 25.625 msec with a frame rate of 100 Hz. Forty MFCC's were extracted from each frame. Thirteen linearly spaced and twenty seven log spaced triangular filters were used to group the FFT bins. The lowest frequency was chosen to be 133.33 Hz, a linear spacing of 66.66 Hz and log spacing of 1.07 were used. In extracting the features we followed the Sphinx III specifications [8]. For the band-pass filter implementation of MFCC, forty fourth-order bandpass-filters (spanning the same frequency range as MFCCs) were used. The BPFs are approximately one-seventh octave, constant Q filters. The filters had to be chosen to be approximately one-seventh octave to match the number of triangular filters used for the standard MFCC features. In the speech versus non-speech classification experiments the first thirteen coefficients were used to perform the classification. A Gaussian mixture model based classifier was used to predict the log-likelihood of each frame belonging to a particular class. The log-likelihoods of all frames in

a segment belonging to the two classes were added to make the final decision. Features from each one second segment were mean subtracted and variance normalized [9].

For the speech recognition task, MFCCs were extracted using code based on the HTK toolkit frontend [10], 23 channels were used. Thirteen MFCC coefficients (including the zeroth coefficient) were mean and variance normalized [9] and delta and acceleration features were computed to form a 39-dimensional feature vector. The BPF-based features were also extracted in a similar way. Thirty-two one-sixth octave filters were used for the filter-bank implementation.

3. Amplitude Compression

Previous work [4], [5], [6] shows that root compression is better than logarithmic compression for noise robustness. In this section we revalidate this result and try to explain the effect of different amplitude compressions in terms of a discrimination measure. Figure 1 shows the performance of the classifier for the speech versus non-speech classification task using MFCC features with root and log compression. We see that root compression is much more robust to noise as compared to log compression. The log function gives large negative values for inputs close to zero and this leads to a spreading of the energy in the transform domain (after DCT).

A simple experiment was devised to show that root compression followed by DCT leads to better compaction of energy. The envelope of a speech segment was amplitude compressed and transformed using DCT. Varying number of transformed coefficients were used to reconstruct the amplitude compressed signal and the reconstruction error was calculated. Figure 2 shows the plot of reconstruction error versus the number of coefficients used for the reconstruction, for both log and root compression. It is clear that root compression followed by DCT leads to better compaction of energy since the reconstruction error using fewer coefficients is much lower as compared to the log case.

Another interesting aspect of amplitude compression is the trade-off between performance in clean conditions and robustness to noise. The trade-off can be better understood in terms of between-class and within-class scatter. We define between-class scatter as the distance between the mean of the two clusters and within-class scatter as the mean of the distances of each data point from the mean of its cluster. The ratio of between-class scatter and within-class scatter is used as a measure of discrimination ability. In clean conditions, more compression leads to lower within-class scatter, the between-class scatter is also lower but since there is no noise to confuse the classifier the accuracy is high. In noisy conditions more compression leads to more errors due to the reduced between-class scatter (the reduction in within-class scatter is not able to offset the reduction in between-class scatter). Table 1 shows the effect of log and different degrees of root compression on the discrimination ability of the MFCC features. As is clear from the table, more compression leads to better performance in clean but this adversely affects the performance in noisy conditions. In the sections that follow, root compression refers to the use of a compression factor, $\alpha = 0.3$, unless stated otherwise.

4. Aliasing and Smoothing

In most audio feature extraction processes the number of samples used to represent each frame is small compared to the original sampled waveform. Given that there will be some loss of information in building a compact representation of the audio signal, the key to

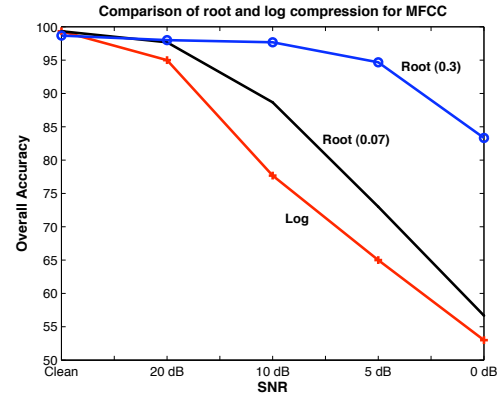


Figure 1: Figure showing the performance of the classifier using MFCC features with log and root compression for a speech versus non-speech classification task. Different SNRs were generated by synthetically adding pink noise.

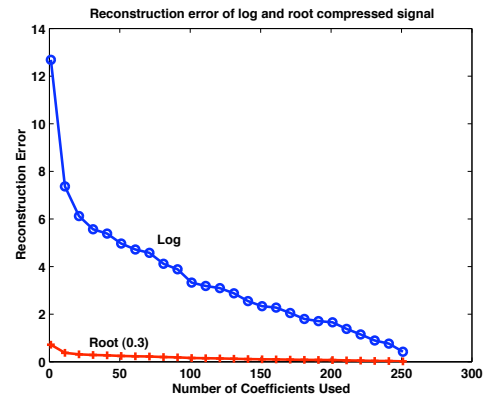


Figure 2: Figure showing that root compression followed by DCT leads to better compaction of energy. Reconstruction error is plotted as a function of number of coefficients used for the reconstruction.

generating better representations is to discard information that is least significant. In case of MFCCs, the FFT and triangular filters lead to discarding of information that is not exactly quantifiable. We also know that due to the sharp peaks of the triangular filters MFCCs are sensitive to small changes in the frequency [11]. The energy estimation in each channel is smoother if frequency decomposition is performed using exponentially spaced bandpass filters and the signal strength in each channel is estimated using a peak detector (implemented using a rectifier and a low-pass filter). Features extracted from a smooth representation are not perturbed by perceptually irrelevant variations in the signal. We know that central auditory neurons cannot respond to very fast temporal modulations [12] and hence smoothing over 10–20ms does not discard perceptually relevant features. The high frequency components that are smoothed out are most likely perceptually insignificant. By low-pass filtering the signal and then down sampling, we are discarding information in a more intelligent way. We refer to features obtained in this manner as BPF-MFCC [13], [14].

Measure of discrimination ability for log and various degrees of root compression		
Compression	Clean	Noisy
Log	0.4364	0.2372
Root ($\alpha = 0.07$)	0.4182	0.2387
Root ($\alpha = 0.3$)	0.3645	0.2473
Root ($\alpha = 0.7$)	0.3196	0.2719

Table 1: Table showing that more compression yields greater discrimination in clean conditions. However, in noisy conditions less compression yields better class discrimination.

Figure 3 shows the improvement in performance of MFCC for the speech versus non-speech classification task by the use of band-pass filters. The BPF-MFCC representation degrades more gracefully with falling SNR.

5. Spatial Derivative

In the previous section, we showed that the performance of BPF-MFCC in noisy conditions is much better than MFCC but the performance in clean conditions is slightly worse. The performance in clean conditions can be improved by introducing yet another processing stage which is directly motivated by physiological processing. The BPFs used for the frequency decomposition are not very sharp and result in some amount of frequency spreading across the channels. This spreading can be limited by the use of a spatial derivative (implemented as a difference operation between adjacent channels). The spatial derivative is used to model the lateral inhibitory network in the cochlear nucleus [12]. We refer to BPF-MFCC with spatial derivative as noise-robust auditory features (NRAF). Apart from limiting the frequency spreading by sharpening the filter response, the spatial derivative stage, in clean and high SNR conditions, enhances the contrast across the spectral profile. This can be thought of as an edge detection operation common in image processing, although the effect in audio is less dramatic due to lack of abrupt changes across frequency channels.

The comparison of MFCC, BPF-MFCC and NRAF is shown in Fig. 3. As predicted the performance of NRAF in clean and moderate SNR cases is better than that of BPF-MFCC, but the performance in high noise case (0 dB SNR) is lower. This can be explained as follows, in very low SNR cases where the noise variance is equal to or greater than the signal variance, the spatial derivative results in some loss of signal component due to subtraction, i.e. the difference operation removes some signal component from channels whose adjacent higher channels are noisy. However if the noise is Gaussian, NRAF performs as well as BPF-MFCC even in low SNR cases.

6. Information-Theoretic Clustering Validity Measure

In this section we use an information theoretic measure of clustering to substantiate the fact that NRAFs are better than the original MFCCs not only in terms of noise-robustness but also in terms of class discrimination ability. Conditional entropy is used as a criterion for evaluating the clustering validity of clustering algorithms [15]. By using a very “naive” clustering algorithm the clustering properties of the underlying attributes can be studied. Mahalanobis distance from the mean of the two clusters is used as the clustering algorithm to study the effect of synthetically added

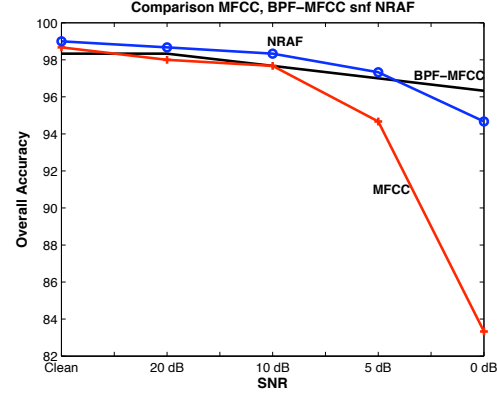


Figure 3: Figure showing the comparative performance of MFCC, BPF-MFCC and NRAF for the speech versus non-speech classification task. Different SNRs were obtained by synthetically adding pink noise. Root compression was used for all the features.

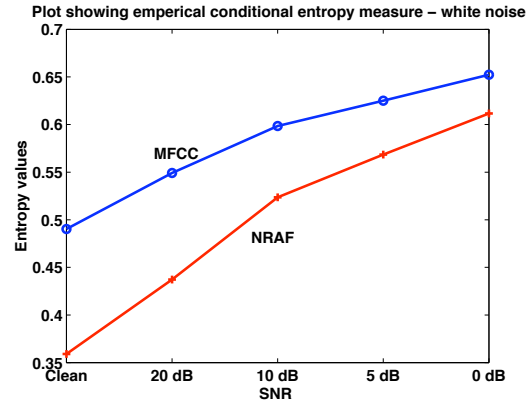


Figure 4: Figure showing the empirical conditional entropy measures for MFCC and NRAF for a two-class, two-cluster case. It is seen that NRAFs outperform MFCCs in the speech versus non-speech classification task. White noise was synthetically added to generate different SNRs.

noise on the clustering properties of MFCC and NRAF.

Given a set of class labels $c \in C$ and clusters $k \in K$, conditional entropy, $H(C|K)$ is approximated by empirical conditional entropy, $H^e(C|K)$ given by,

$$H^e(C|K) = - \sum_{c=1}^{|C|} \sum_{k=1}^{|K|} \frac{h(c,k)}{n} \log \frac{h(c,k)}{h(k)}$$

where, $h(c,k)$ is the number of examples of class c assigned to cluster k , n is the total number of examples and $h(k)$ is the associated marginal. Conditional entropy gives the number of bits required to encode all the class information given we know the clusters $\{k_i\}$ and the model $\Pi = \{h(c,k)\}$. A lower value of conditional entropy indicates that the cluster labels are good indicators of the class labels. For the case of two classes and two clusters, the empirical conditional entropy measure is shown in Fig. 4. As is evident, NRAF has a lower value of conditional en-

Recognition results on Aurora 2 (clean condition)								
	Set A		Set B		Set C		Average	
	MFCC	NRAF	MFCC	NRAF	MFCC	NRAF	MFCC	NRAF
Clean	98.70 %	98.38 %	98.79 %	98.38 %	98.72 %	98.51 %	98.74 %	98.42 %
20 dB	95.88 %	94.86 %	96.08 %	94.88 %	96.04 %	94.61 %	96.00 %	94.78 %
15 dB	92.58 %	90.76 %	92.99 %	90.76 %	93.13 %	90.96 %	92.89 %	90.82 %
10 dB	84.53 %	83.73 %	85.23 %	82.58 %	85.54 %	83.95 %	85.10 %	83.41 %
5 dB	68.64 %	68.46 %	68.24 %	65.67 %	70.33 %	70.61 %	69.07 %	68.25 %
0 dB	41.82 %	40.03 %	42.59 %	36.75 %	43.18 %	43.93 %	42.53 %	40.24 %
-5 dB	17.31 %	16.19 %	17.55 %	13.41 %	19.71 %	17.96 %	18.19 %	15.85 %

Table 2: Two Gaussian components per mixture was used for every state, except silence, which was modelled using 4 components. Only 75 % of the training data was utilized, but the whole test set was used for evaluation. It is clear that a low complexity backend is able to fit the MFCC data better than NRAF data.

Recognition results on Aurora 2 (clean condition)								
	Set A		Set B		Set C		Average	
	MFCC	NRAF	MFCC	NRAF	MFCC	NRAF	MFCC	NRAF
Clean	99.06 %	99.05 %	99.06 %	99.05 %	99.16 %	99.10 %	99.09 %	99.06 %
20 dB	96.66 %	96.45 %	97.14 %	96.55 %	96.77 %	96.32 %	96.86 %	96.44 %
15 dB	93.64 %	93.62 %	94.40 %	93.28 %	94.06 %	93.30 %	94.03 %	93.40 %
10 dB	86.17 %	87.32 %	87.76 %	86.64 %	87.28 %	87.86 %	87.07 %	87.27 %
5 dB	71.23 %	73.52 %	72.36 %	71.48 %	73.08 %	75.62 %	73.08 %	73.54 %
0 dB	44.84 %	48.16 %	45.52 %	46.23 %	46.36 %	52.41 %	45.57 %	48.93 %
-5 dB	19.28 %	21.85 %	19.40 %	20.37 %	20.42 %	23.89 %	19.70 %	22.04 %

Table 3: Three Gaussian components per mixture was used for every state, except silence, which was modelled using 6 components. The entire training and test data was used. It is seen that when the complexity of the backend is increased it is able to model the NRAF data better.

trophy, implying that it clusters better than MFCCs in clean and noisy conditions.

7. Speech Recognition Results

The speech recognition results for the Aurora 2 task in clean training condition are presented in Tables 2-4 below. MFCC and NRAF features were MVA processed as suggested by Chen et al. [9]. Delta and acceleration coefficients were extracted from the MVA processed static features. The zeroth coefficient was used since it is shown to respond better to MVA than using the log energy. Logarithmic compression was used for both feature sets.

In order to evaluate the performances of MFCC and NRAF, we first trained a HMM with 2 Gaussian components per mixture for every state and 4 components for the silence model. The HMM was trained on only 75 % of the training data but tested on the whole test set. Next, we trained HMMs with 3/6 components per mixture for every state and 6/12 components for the silence model respectively. These HMMs were trained on the complete training set and evaluated on the whole test set. We see that MFCC's outperform NRAF's in the first case, but as the complexity of the backend increases NRAF outperforms MFCC as seen in Fig. 5. The increased modeling ability of the backend enables the recognizer to better fit the extra information encoded by the NRAF representation. We hypothesize that detailed auditory model representations, in general, fail to beat MFCC's performance on speech recognition tasks due to the fact that these representation owing to their encoding of more information than MFCC's need more com-

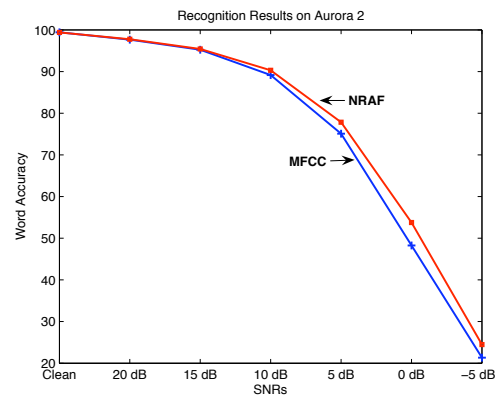


Figure 5: Figure showing the performance of NRAF and MFCC on Aurora 2 task. Six Gaussian mixture was used for each state and silence was modelled using 12 component mixture.

plexity in the backend to learn the data better. This directly implies the need for more training data in order to avoid overfitting.

8. Conclusions and Future Work

This paper justifies three improvements to MFCC features that are motivated by a more accurate auditory model. We tested these

Recognition results on Aurora 2 (clean condition)									
	Set A		Set B		Set C		Average		Relative
	MFCC	NRAF	MFCC	NRAF	MFCC	NRAF	MFCC	NRAF	Improvement
Clean	99.41 %	99.36 %	99.41 %	99.36 %	99.44 %	99.42 %	99.42 %	99.38 %	-0.04 %
20 dB	97.53 %	97.72 %	97.89 %	97.79 %	97.49 %	97.78 %	97.64 %	97.76 %	0.12 %
15 dB	95.02 %	95.45 %	95.47 %	95.44 %	95.19 %	95.37 %	95.23 %	95.42 %	0.20 %
10 dB	88.55 %	90.16 %	89.83 %	89.80 %	89.15 %	90.90 %	89.17 %	90.29 %	1.26 %
5 dB	74.47 %	77.75 %	75.18 %	76.60 %	75.63 %	79.20 %	75.09 %	77.85 %	3.68 %
0 dB	47.94%	53.89 %	48.28 %	51.02 %	48.55 %	56.41 %	48.26 %	53.77 %	11.42 %
-5 dB	20.94 %	24.56 %	21.10 %	22.33 %	22.01 %	26.51 %	21.35 %	24.47 %	14.61 %

Table 4: Six Gaussian components per mixture was used for every state, except silence, which was modelled using 12 components. The entire training and test data was used. The increased modeling ability of the backend enables it to better fit the extra information encoded by the NRAF representation.

changes using a simple speech versus non-speech test and a speech recognition task. Replacing the log compression with a root compression improves the noise-robustness of MFCCs. Low-pass filtering the signal in each channel before decimation avoids aliasing and leads to a smoother signal envelope in each channel. Further, the benefit of spatial derivative in clean and high to moderate SNR cases is demonstrated. Future work would involve developing BPF-based implementation of MFCC with different time constants for the LPF in different channels, hopefully leading to an even better representation.

9. References

- [1] Nima Mesgarani, Shihab Shamma, and Malcolm Slaney, "Speech discrimination based on multiscale spectro-temporal modulations," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, Montreal, Canada, May 2004.
- [2] Mark D. Skowronski and John G. Harris, "Exploiting independent filter bandwidth of human factor cepstral coefficients in automatic speech recognition," *Journal of Acoustical Society of America*, vol. 116, no. 3, pp. 1774–1780, sept 2004.
- [3] M.J. Hunt, M. Lenning, and P. Mermelstein, "Experiments in syllable-based recognition of continuous speech," in *IEEE International Conference on Accoustics, Speech, and Signal Processing*, Denver, CO, apr 1980.
- [4] J.S. Lim, "Spectral root homomorphic deconvolution system," *IEEE Trans. ASSP*, vol. 27, no. 3, pp. 223–233, 1979.
- [5] P. Alexandre and P. Lockwood, "Root cepstral analysis: A unified view," *Speech Communication*, vol. 3, pp. 277–288, 1993.
- [6] Ruhi Sarikaya and John H.L. Hansen, "Analysis of the root-cepstrum for acoustic modeling and fast decoding in speech recognition," in *Eurospeech*, Aalborg, Denmark, sept 2001.
- [7] Masataka Goto, Hiroki Hashiguchi, Takuichi Nishimura, and Ryuichi Oka, "RWC music database: Music genre database and musical instrument sound database," in *Proceedings of the 4th International Conference on Music Information Retrieval*, oct 2003, pp. 229–230.
- [8] Michael Seltzer, "Sphinx III signal processing front end specification http://cmusphinx.sourceforge.net/sphinx3/s3_fe_spec.pdf," .
- [9] C.-P Chen, K. Filali, and J. A. Bilmes, "Frontend post-processing and backend model enhancement on the aurora 2.0/3.0 databases," in *International Conference on Speech and Language Processing*, 2002, pp. 241–244.
- [10] ETSI ES 201 108 v1.1.3 (2003-09), "Speech Processing, Transmission and Quality Aspects (STQ); Distributed speech recognition; Front-end feature extraction algorithm; Compression algorithms," .
- [11] Steven Beet, "Email contribution to auditory mailing list, Nov, 2004. <http://www.auditory.org/postings/2004/833.html>," .
- [12] Xiaowei Yang, Kuansan Wang, and Shihab Shamma, "Auditory representations of acoustic signals," *IEEE Transactions on Information Theory*, vol. 38, no. 2, pp. 824–839, March 1992.
- [13] Sourabh Ravindran, Cenk Demiroglu, and David Anderson, "Speech recognition using filter-bank features," in *Asilomar Conference on Signals and Systems*, Pacific Grove, CA, Nov. 2003.
- [14] Paul Smith, Matt Kucic, Rich Ellis, Paul Hasler, and David V. Anderson, "Cepstrum frequency encoding in analog floating-gate circuitry," in *Proceedings of the IEEE International Symposium on Circuits and Systems*, Phoenix, AZ, may 2002, vol. IV, pp. 671–674.
- [15] Byron E. Dom, "An information-theoretic external cluster-validity measure," in *IBM Research Report RJ 10219*, May 2001.