



Novel Subband Autoencoder Features for Non-intrusive Quality Assessment of Noise Suppressed Speech

Meet H. Soni, Hemant A. Patil

Dhirubhai Ambani Institute of Information and Communication Technology, India

{meet_soni, hemant_patil}@daiict.ac.in

Abstract

In this paper, we propose a novel feature extraction architecture of Deep Neural Network (DNN), namely, subband autoencoder (SBAE). The proposed architecture is inspired by the Human Auditory System (HAS) and extracts features from speech spectrum in an unsupervised manner. We have used features extracted by this architecture for non-intrusive objective quality assessment of noise suppressed speech signal. The quality assessment problem is posed as a *regression* problem in which mapping between the acoustic features of speech signal and the corresponding subjective score is found using single layer Artificial Neural Network (ANN). We have shown experimentally that proposed features give more powerful mapping than Mel filterbank energies, which are state-of-the-art acoustic features for various speech technology applications. Moreover, proposed method gives more accurate and correlated objective scores than current standard objective quality assessment metric ITU-T P.563. Experiments performed on NOIZEUS database for different test conditions also suggest that objective scores predicted using proposed method are more robust to different amount and types of noise.

Index Terms: speech quality assessment, autoencoder, regression.

1. Introduction

Speech quality assessment is very important in many applications including telephone networks, Voice over Internet Protocol (VoIP), multimedia applications, etc. The best way to assess the quality of speech is to take the opinion of human listeners. To do this task, listening tests are conducted which serves as a subjective quality assessment measure. The widely used subjective measure is Mean Opinion Score (MOS). However, some fundamental difficulties including cost, time consumption and in some cases, the reliability of the subjective test (due to cognitive factors associated with the listener), makes it unsuitable for several applications which require in-service, real-time or in-process quality assessment. Hence, to overcome these limitations, there is a requirement for a reliable objective measure to assess the speech quality. Objective speech quality assessment has attracted researchers over past two decades [1–7].

The aim of objective quality evaluation is to find the replacement for human judgment of perceived speech quality. Objective evaluation techniques are less complex, less expensive in terms of resources and time complexity and give more consistent results [8]. Objective evaluation techniques are categorized in two ways, namely, intrusive and non-intrusive. Intrusive assessments are based on waveform comparison wherein reference speech signal is available for comparison. On the

other hand, non-intrusive quality assessment (also known as single-ended, no-reference or output-based quality assessment) is performed using single speech waveform, without any reference or the ground truth. Intrusive methods are more straightforward, less complex and relatively more accurate than the non-intrusive ones. However, in many practical scenarios such as wireless communication, Voice over IP (VoIP) and other in-service applications (requiring monitoring of speech quality), intrusive methods cannot be applied due to unavailability of reference speech signal. In such realistic scenarios, it is necessary to have a reliable non-intrusive method for quality assessment.

An early attempt towards non-intrusive assessment of speech based on spectrogram analysis is presented in [1]. The study reported in [2] uses Gaussian Mixture Models (GMMs) to create artificial reference model to compare degraded speech signals whereas in [3], speech quality is being predicted by Bayesian inference and minimum mean square estimation (MMSE) based on trained GMMs. In [4], a speech quality assessment algorithm based on temporal envelope representation of speech is presented. Different features extracted from speech have been detected to be useful for speech quality assessment. Spectral dynamics, spectral flatness, spectral centroid, spectral variance, pitch and excitation variance was used for quality prediction in [5]. The authors in [9] used perceptual linear prediction (PLP) coefficients for quality assessment. Quality assessment problem is posed as a regression problem and the mapping between acoustic features and subjective score was found in [10] using MFCCs and in [8] using filterbank energies as acoustic feature. To find mapping, Support Vector Regression (SVR) was used. Bag of Words (BoW) inspired codebook approach was presented in [11]. While authors in [12] used spectro-temporal features for the same task. Several combinations of auditory features was used in [13] for the quality assessment task.

Recently, deep learning methods are gaining popularity for feature extraction from the raw data. Autoencoder is such network which uses Deep Neural Network (DNN) or Restricted Boltzmann Machine (RBM) to extract low-dimensional information from high-dimensional raw data [14–17]. Autoencoder has been widely used for automatic speech recognition (ASR) systems for noisy or reverberant conditions. In [18] and [19], authors used autoencoder as denoising frontend for such ASR task. Autoencoder was used to find a mapping between the spectrum of noisy speech and clean speech for noise reduction in ASR system [20]. Autoencoder was also used for speech enhancement task in [21]. For speech coding, autoencoder was used to encode speech spectrum in [22]. Very recently, the study reported in [23] showed the use of autoencoder for noise reduction in speaker verification system. Deep autoencoder was

used in [24] for noise aware training for ASR in the noisy environment. Features learned by deep autoencoder were used for Statistical Parametric Speech Synthesis (SPSS) using DNN in [25]. Despite these properties and their usefulness, autoencoder features are seldom used as primary acoustic features for any speech technology application. Primary reason for this unpopularity is the manner in which autoencoder extracts the features from speech spectrum. Features extracted by autoencoder are difficult to interpret since there is no control over its learning.

To overcome these limitations, many variants of the autoencoder is proposed to use prior knowledge of speech-domain. A new architecture called transforming autoencoder was used in [26] to detect acoustic events in speech signal for ASR task. Phone recognition task was done using mean-covariance RBM in [27]. The study reported in [28] proposed architecture of autoencoder in which decoding block was constrained for stretching and compressing frequency-domain for ASR task.

In this paper, we propose a new architecture called *subband autoencoder*, which is closer to the Human Auditory System (HAS). In proposed architecture, we have constrained the connectivity of the units of input layer to the units of the first hidden layer of autoencoder. By doing this, each unit in the first hidden layer is forced to capture information about a particular band of the speech spectrum which mimics human auditory processing in some sense. We have used features extracted by proposed architecture for the quality assessment task. The problem of speech quality assessment is posed as a regression problem, same as previously done in [10] and [8]. However, we have used proposed features as the acoustic features and used a single-layer artificial neural network (ANN) as a regression model. ANN was chosen due to its universal approximation abilities and need of least tuning of the parameters. We have shown experimentally that proposed features provide more variability in feature vectors for speech signals having different types and amount of noise. Moreover, they are able to reconstruct speech spectrum more precisely than filterbank energies. These properties of subband autoencoder features suggest that they capture noise information in a better way than MFCCs or filterbank energies.

2. Proposed Subband Autoencoder

2.1. Architecture of subband autoencoder

Fig. 1 shows the architecture of proposed subband autoencoder. The main difference between proposed architecture and architecture of an autoencoder [15] is the connectivity of neurons or units immediately after the input layer. In autoencoder, each unit in the layer immediately after input layer is connected with all the units of the previous layer. While in the case of proposed subband autoencoder, the connectivity is restricted. In proposed architecture, each unit of the first hidden layer is connected with a particular frequency band of the input spectrum. Hence, each unit in the first layer will encode the information about that particular frequency band only, with which it is connected. The decoding structure is same as a general autoencoder with full connectivity. The band structure of *restricted* connectivity for neurons is same as Mel filterbank, implying one neuron in the first layer is connected with the frequencies of one Mel filterbank. This architecture is nearer to HAS and provides more meaningful information than autoencoder in the case of speech. Mathematically, operation of the subband layer can be represented as follows:

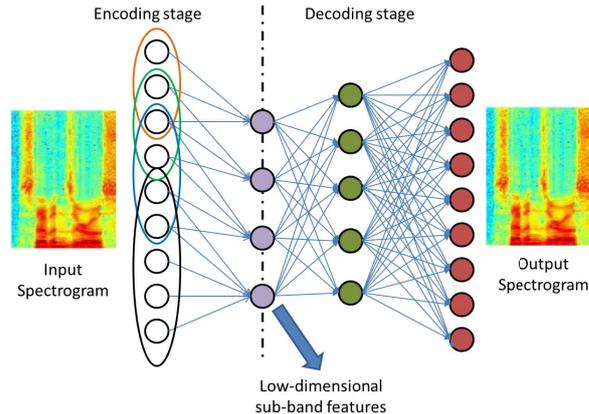


Figure 1: Proposed architecture of subband autoencoder.

$$a_i = f\left(\sum_j W_{ij}^1 \times x_j\right), \quad (1)$$

where a_i is i^{th} subband feature, x_j is short-time power corresponding to j^{th} filterbank frequencies and W_{ij}^1 are weights corresponding to i^{th} subband feature. f represents nonlinear activation function of the neuron. The functionality of preceding layers of subband autoencoder is same as that of a simple autoencoder [15]. Proposed architecture can be trained by back-propagation similarly to an autoencoder. a_j learned by subband autoencoder can be used as low-dimensional features for speech technology task. These features are different from filterbank energies in following ways: The first difference is in the method of extracting features. MFCCs or filterbank energies are hand-crafted features while subband autoencoder features are learned by a machine learning approach. Filterbank energies are extracted in a linear way, while subband autoencoder features are extracted in a nonlinear manner. The latter property may provide some more useful information about speech spectrum variations for different conditions, such as speech signals having different types and amount of noise.

2.2. Analysis of subband autoencoder features

Fig. 2 shows the subband autoencoder features and filterbank energies for clean and noisy speech along with short-time spectrum. Fig. 3 shows mean of subband autoencoder features and filterbank energies of 20 clean and noisy speech utterances having different Signal-to-Noise Ratio (SNR). As it can be observed, both the features vary under the influence of additive noise. Hence, both of them can be used for measuring the quality of speech signal they represent. However, it is difficult to deduce which one of them will be better suitable to represent overall quality of speech just by observing the features. In [10], authors have justified the use of MFCCs for quality assessment task by arguing about their ability to reconstruct speech spectrum [29, 30]. MFCCs are able to restore the perceptual quality of underlying clean speech from given noisy speech. The same argument can be made for filterbank energies. Proposed subband autoencoder features are also invertible and hence, they can be also used quality assessment task. Fig. 2 shows original and inverted short-time spectrum using proposed features and filterbank energies for clean as well as noisy speech. By visual inspection, it can be observed that spectrum inverted using proposed features is more identical to the original spectrum for clean as well as noisy speech. To quantify the similarity

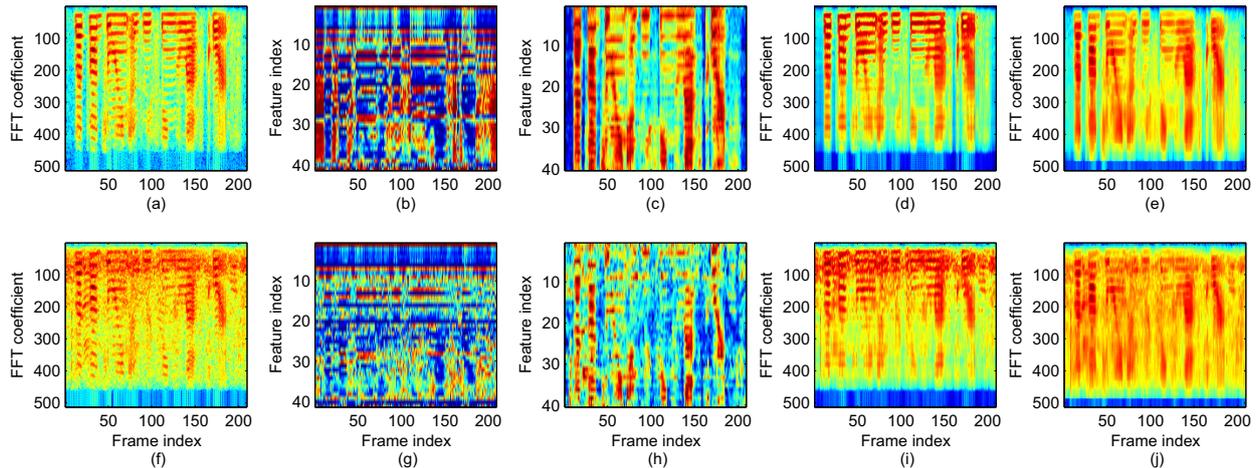


Figure 2: (a) Short-time spectrum, (b) features extracted using proposed subband autoencoder and (c) mel-filterbank energies for clean speech. Reconstructed short-time spectrum using (d) SBAE features and (e) filterbank energies for clean speech. Similarly, (f) short-time spectrum, (g) subband autoencoder features and (h) mel-filterbank energies and reconstructed short-time spectrum using (d) SBAE features and (e) filterbank energies for noisy speech, corrupted with additive car noise of 5 dB SNR.

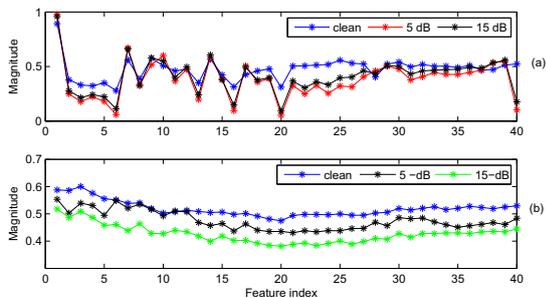


Figure 3: Mean (a) subband autoencoder features and (b) filterbank energies of 20 speech utterances having different amount of additive car noise.

between original and inverted spectrum using both the features, we have calculated log-spectral distortion (LSD) between two spectra. LSD between original spectrum and inverted spectrum of clean speech was 0.68 dB using proposed features, while it was 0.87 dB in case of filterbank energies. In case of noisy speech it was 0.47 dB and 1.01 dB using proposed features and filterbank energies, respectively. Each spectrum was normalized between $0-1$ to make dynamic range uniform. It shows that proposed subband autoencoder features are able to invert speech spectrum more precisely than filterbank energies, and for that matter, MFCCs. Hence, it can be said that they capture the underlying information of speech spectrum in a better way. Similar studies were shown for traditional autoencoder features in very recent work [25].

3. Experimental Results

3.1. Experimental setup

All experiments were performed on NOIZEUS database [6]. The database had speech utterances which were corrupted by different kind and different amount of additive noise. It also had speech utterances which were enhanced by different noise suppression algorithms. The speech utterances were corrupted

by four types of additive noise with two SNR levels. The noise suppression algorithms fall under four different classes. A complete description of database can be found in [8] and noise suppression algorithms in [6], [31]. The subjective evaluation of the speech utterances was performed according to ITU-T Recommendation P.835 [32], [6]. Both subband autoencoder features and filterbank energies were extracted from this database. The performance comparison of the both features was done for $40-D$ (dimensional) features. To extract subband autoencoder features from FFT spectrum, proposed architecture was used. The configuration of subband autoencoder was $513-40-200-513$, meaning 513 units in the first layer, 40 units in the second layer, which is subband layer and so on. All units had sigmoid as nonlinearity in all the layers. To demonstrate the ability of proposed features to capture general spectral information, subband autoencoder was trained only using 150 utterances which were not used for further experiments. Mel filterbank energies of same dimensions were extracted from speech utterances. Moreover, for comparison purpose, a simple autoencoder is also used to extract $40-D$ features. The architecture of the autoencoder was $513-250-40-250-513$, meaning 513 units in input layer, 250 units in first hidden layer, 40 units in second hidden layer (to extract $40-D$ autoencoder features) and so on.

To find a mapping between the mean of features extracted from speech and their subjective score, artificial neural network (ANN) with single hidden layer was used. ANN was used due to its universal approximation strength. The mean of the features over different conditions was used to find mapping. A number of hidden units in ANN was 350 which was selected using validation data. The network was regularized using standard weight decay method to prevent overfitting. Per condition speech signals and their subjective scores were considered for experiments as used in [8]. Total 4 different tests were performed to check the robustness of proposed algorithm. First test was standard 8-fold cross-validation [8]. Other 3 tests were also performed with similar conditions as shown in [8].

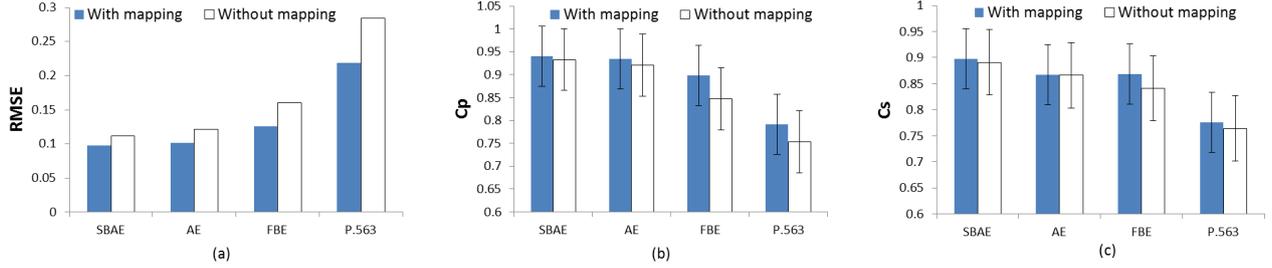


Figure 4: (a) RMSE, (b) C_p and (c) C_s between predicted scores and actual subjective scores using proposed features (SBAE), autoencoder features (AE), filterbank energies (FBE) and ITU-T P.563 standard. Results are shown with and without polynomial mapping. C_p and C_s are shown with 95 % confidence intervals.

Table 1: RMSE, C_p and C_s for test 2, 3 and 4. Results are shown with and without polynomial mapping.

Test 2 Without mapping				Test 2 With mapping			
Method	RMSE	C_p	C_s	Method	RMSE	C_p	C_s
SBAE	0.12	0.92	0.92	SBAE	0.11	0.94	0.92
AE	0.13	0.92	0.91	AE	0.11	0.93	0.91
FBE	0.16	0.84	0.85	FBE	0.13	0.90	0.82
P563	0.37	0.66	0.65	P563	0.29	0.66	0.65
Test 3 Without mapping				Test 3 With mapping			
Method	RMSE	C_p	C_s	Method	RMSE	C_p	C_s
SBAE	0.23	0.82	0.84	SBAE	0.18	0.88	0.87
AE	0.26	0.73	0.72	AE	0.20	0.82	0.76
FBE	0.27	0.69	0.70	FBE	0.24	0.72	0.69
P563	0.36	0.68	0.68	P563	0.33	0.68	0.68
Test 4 Without mapping				Test 4 With mapping			
Method	RMSE	C_p	C_s	Method	RMSE	C_p	C_s
SBAE	0.18	0.86	0.83	SBAE	0.16	0.86	0.83
AE	0.23	0.76	0.79	AE	0.21	0.78	0.79
FBE	0.23	0.75	0.78	FBE	0.21	0.78	0.78
P563	0.32	0.68	0.65	P563	0.30	0.69	0.65

3.2. Results and discussions

To evaluate the performance, three common criteria were used: Pearson linear correlation coefficient C_p (for prediction accuracy), Spearman rank order correlation coefficient C_s (for prediction monotonicity) and Root Mean Squared Error (RMSE) between predicted objective score and subjective scores [6]. For an ideal match between the objective and subjective scores, $C_p=C_s=1$ and RMSE=0. Moreover, results are also shown with 3^{rd} order polynomial mapping suggested in [7] to eliminate offset between subjective and objective scores.

Fig. 4 shows RMSE, C_p and C_s calculated for subband autoencoder features, autoencoder features and filterbank energies for test 1. We also compared our results with ITU P.563 standard [7], which is the standard objective measure for non-intrusive speech quality assessment. C_p and C_s are shown with 95 % confidence intervals. Objective scores predicted using proposed features are more accurate as well as more correlated with actual subjective scores. Moreover, the overlap of 95 % confidence intervals between P.563 and proposed metric is almost zero. Hence, it can be said that proposed metric is nearer to actual subjective scores than state-of-the-art P.563 score, while it is not the case using filterbank energies. Table 1 shows RMSE, C_p and C_s calculated for test conditions 2, 3 and 4. Results of different test conditions suggest that both P.563 metric and scores predicted using filterbank energies as well as using proposed features are condition-dependent. While the performance of each the metric varies according to conditions, proposed metric is more robust to different noisy condi-

tions than the other three. Hence, objective scores predicted by proposed method are more reliable as compared to the other methods. These results are in coherence with the analysis of proposed features presented in Section 2. Moreover, features of proposed architecture of autoencoder perform slightly better than a traditional autoencoder for *test 1* and *test 2*. While in *test 3* and *test 4*, proposed features perform significantly better than traditional autoencoder features. This suggests that proposed architecture is more useful than unconstrained autoencoder for the task. The reason behind this can be the constraint on autoencoder to learn band wise information from input spectrum.

4. Summary and Conclusions

In this paper, we have proposed new feature extraction architecture of DNN, which is inspired by the Human Auditory System (HAS). We have used features extracted using proposed architecture for the non-intrusive objective quality assessment task. Here, quality assessment problem is posed as a regression problem, and proposed features are used as acoustic features to find a mapping between speech signal and its subjective quality score. Results of our experiments show that proposed features give more powerful mapping than state-of-the-art acoustic features. Moreover, proposed metric also gives more accurate and correlated objective score than current baseline non-intrusive objective metric ITU-T P.563. Our future work includes testing the performance of proposed features for other intrusive and non-intrusive quality assessment task such as quality assessment of vocoder, speech synthesis system, etc.

5. References

- [1] O. Au and K. Lam, "A novel output-based objective speech quality measure for wireless communication," in *Fourth International Conference on Signal Processing Proceedings (ICSP)*, Beijing, China, 1998, pp. 666–669.
- [2] T. H. Falk, Q. Xu, and W.-Y. Chan, "Non-intrusive gmm-based speech quality measurement," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Philadelphia, Pennsylvania, USA, 2005, pp. 125–128.
- [3] G. Chen and V. Parsa, "Bayesian model based non-intrusive speech quality evaluation," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Philadelphia, Pennsylvania, USA, 2005, pp. 385–388.
- [4] D.-S. Kim, "Anique: An auditory model for single-ended speech quality estimation," *IEEE Transactions on Speech and Audio Processing*, vol. 13, no. 5, pp. 821–831, 2005.
- [5] V. Grancharov, D. Y. Zhao, J. Lindblom, and W. B. Kleijn, "Low-complexity, nonintrusive speech quality assessment," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 6, pp. 1948–1956, 2006.
- [6] Y. Hu and P. C. Loizou, "Evaluation of objective quality measures for speech enhancement," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 1, pp. 229–238, 2008.
- [7] T. Falk and W. Chan, "Single ended method for objective speech quality assessment in narrowband telephony applications," *ITU-T*, 2004.
- [8] M. Narwaria, W. Lin, I. V. McLoughlin, S. Emmanuel, and L.-T. Chia, "Nonintrusive quality assessment of noise suppressed speech with mel-filtered energies and support vector regression," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 4, pp. 1217–1232, 2012.
- [9] T. H. Falk and W.-Y. Chan, "Single-ended speech quality measurement using machine learning methods," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 6, pp. 1935–1947, 2006.
- [10] M. Narwaria, W. Lin, I. V. McLoughlin, S. Emmanuel, and C. L. Tien, "Non-intrusive speech quality assessment with support vector regression," in *Advances in Multimedia Modeling*, 2010, pp. 325–335.
- [11] Q. Li, W. Lin, Y. Fang, and D. Thalmann, "Bag-of-words representation for non-intrusive speech quality assessment," in *2015 IEEE China Summit and International Conference on Signal and Information Processing (ChinaSIP)*, Chengdu, China: IEEE, 2015, pp. 616–619.
- [12] Q. Li, Y. Fang, W. Lin, and D. Thalmann, "Non-intrusive quality assessment for enhanced speech signals based on spectro-temporal features," in *2014 IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*, IEEE, 2014, pp. 1–6.
- [13] R. K. Dubey and A. Kumar, "Non-intrusive speech quality assessment using several combinations of auditory features," *International Journal of Speech Technology (IJST)*, vol. 16, no. 1, pp. 89–101, 2013.
- [14] J. Gehring, Y. Miao, F. Metze, and A. Waibel, "Extracting deep bottleneck features using stacked auto-encoders," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Vancouver, Canada, 2013, pp. 3377–3381.
- [15] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [16] T. N. Sainath, B. Kingsbury, and B. Ramabhadran, "Auto-encoder bottleneck features using deep belief networks," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Kyoto, Japan, 2012, pp. 4153–4156.
- [17] D. Yu and M. L. Seltzer, "Improved bottleneck features using pre-trained deep neural networks," in *INTERSPEECH*, 2011, pp. 237–240.
- [18] T. Ishii, H. Komiya, T. Shinozaki, Y. Horiuchi, and S. Kuroiwa, "Reverberant speech recognition based on denoising autoencoder," in *INTERSPEECH*, Lyon, France, 2013, pp. 3512–3516.
- [19] X. Feng, Y. Zhang, and J. Glass, "Speech feature denoising and dereverberation via deep autoencoders for noisy reverberant speech recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Florence, Italy, 2014, pp. 1759–1763.
- [20] A. L. Maas, Q. V. Le, T. M. O’Neil, O. Vinyals, P. Nguyen, and A. Y. Ng, "Recurrent neural networks for noise reduction in robust asr," in *INTERSPEECH*, Portland, Oregon, 2012, pp. 22–25.
- [21] X. Lu, Y. Tsao, S. Matsuda, and C. Hori, "Speech enhancement based on deep denoising autoencoder," in *INTERSPEECH*, 2013, pp. 436–440.
- [22] L. Deng, M. L. Seltzer, D. Yu, A. Acero, A.-R. Mohamed, and G. E. Hinton, "Binary coding of speech spectrograms using a deep auto-encoder," in *INTERSPEECH*, Makuhari, Japan, 2010, pp. 1692–1695.
- [23] O. Plchot, L. Burget, H. Aronowitz, and M. Pavel, "Audio enhancing with DNN autoencoder for speaker recognition," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Shanghai, China, 2016, pp. 5090–5094.
- [24] K. H. Lee, S. J. Kang, W. H. Kang, and N. S. Kim, "Two-stage noise aware training using asymmetric deep denoising autoencoder," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Shanghai, China, 2016, pp. 5765–5769.
- [25] S. Takaki and J. Yamagishi, "A deep auto-encoder based low-dimensional feature extraction from FFT spectral envelopes for statistical parametric speech synthesis," in *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Shanghai, China, 2016, pp. 5090–5094.
- [26] N. Jaitly and G. E. Hinton, "A new way to learn acoustic events," *Advances in Neural Information Processing Systems (NIPS)*, vol. 24, 2011.
- [27] G. Dahl, A.-r. Mohamed, G. E. Hinton *et al.*, "Phone recognition with the mean-covariance restricted boltzmann machine," in *Advances in Neural Information Processing Systems (NIPS)*, 2010, pp. 469–477.
- [28] N. Jaitly and G. E. Hinton, "Using an autoencoder with deformable templates to discover features for automated speech recognition," in *INTERSPEECH*, Lyon, France, 2013, pp. 1737–1740.
- [29] L. E. Boucheron and P. L. De Leon, "On the inversion of mel-frequency cepstral coefficients for speech enhancement applications," in *International Conference on Signals and Electronic Systems (ICSES)*, Krakow, Poland, 2008, pp. 485–488.
- [30] X. Shao and B. Milner, "Clean speech reconstruction from noisy mel-frequency cepstral coefficients using a sinusoidal model," in *IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003 (ICASSP)*, vol. 1, Hong Kong, Hong Kong, 2003, pp. 700–704.
- [31] Y. Hu and P. C. Loizou, "Subjective comparison and evaluation of speech enhancement algorithms," *Speech Communication*, vol. 49, no. 7, pp. 588–601, 2007.
- [32] ITU-T, "ITU-T Rec 835, subjective test methodology for evaluating speech communication systems that include noise suppression algorithm," <http://www.itu.int/rec/T-REC-P.835-200311-I>, 2003, {Last Accessed: 30th March, 2016}.