

# WAVELET IN CONJUNCTION WITH NEURAL NETWORK METHOD FOR SPEECH ENHANCEMENT QUALITY EVALUATION

*K. Daqrouq and Ghada Amer*

Philadelphia Univ., Jordan-haleddaq@yahoo.com  
Benha Univ., Egypt, dr\_ghada@benha-univ.edu.eg

## ABSTRACT

Wavelet Neural Network Evaluation method WNNEM is proposed as a powerful tool for enhanced speech signal evaluation. This objective evaluation measure utilizes Feed forward back Propagation Neural Network FFNN to train the free of noise signal, and then enhanced signal is simulated to the training output results taken for given target. The distance between simulation and the target, over different wavelet sub bands is studied. Four known speech enhancement method for studying the performance of WNNEM are utilized. The advantage of this method is the evaluation of different band passes of frequency based on wavelet transform by neural network, which is very influential tool for non stationary signals processing. Several objective measures are used to investigate the WNNEM compatibility. Results proved the validity of the proposed method.

**Index Terms**— *Wavelet, Neural Network, Speech Enhancement, Quality Evaluation*

## 1. INTRODUCTION

The types of deformation introduced by speech enhancement algorithms can be broadly divided into two kinds: the distortions that change the speech signal itself (called speech distortion) and the distortions that change the background noise (called noise distortion). Of these two types of deformation, listeners seem to be influenced the most by the speech distortion when making judgments of overall quality [1], [2]. the most accurate method for evaluating speech quality is through subjective listening tests. Although subjective evaluation of speech enhancement algorithms is often accurate and reliable (i.e., repeatable) provided it is performed under stringiest conditions (e.g., sizeable listener panel, inclusion of anchor conditions, etc. [4]–[7]), it is costly and time consuming. For that reason, much

effort has been placed on developing objective measures that would predict speech quality with high correlation. Many objective speech quality measures have been proposed in the past to predict the subjective quality of speech [4]. Most of these measures, however, were developed for the purpose of evaluating the distortions introduced by speech codecs and/or communication channels [7]–[12].

## 2. QUALITY EVALUATION SPEECH ENHANCEMENT SYSTEMS

Different methods have been proposed for speech enhancement systems evaluation. All of these methods are based on comparison of original signal with enhanced signal by relative ratio measure or distance measure. The most popular measure, which gives a measure of the signal power improvement related to the noise power is *SNR* [13], and segmental *SNR* (segSNR) [14]. From spectral domain evaluation algorithm, we can mention Weighted Slope Spectral distance (WSS) [15]

$$d_{WSS} = \frac{1}{M} \sum_{m=0}^{M-1} \frac{\sum_{l=1}^K W(I, M) (s_c(I, M) - s_p(I, M))^2}{\sum_{l=1}^K W(I, M)} \quad (1)$$

Where  $W(I, M)$  is the weight placed on  $l$ th frequency band,  $K$  is the number of bands and  $M$  is the number of frames in the signal.  $s_c(I, M)$  and  $s_p(I, M)$  spectral are the slope of the clean and enhanced signals, respectively. Hu and Loizou in [3], used the value of  $K$  as 25.

Cepstrum distance has been used as a difference of original signal cepstrum and enhanced signal cepstrum [3]

$$d_{CEP}(\hat{C}_C, \hat{C}_P) = \frac{10}{\log 0} \sqrt{2 \sum_{k=1}^P (C_C(k) - C_P(k))^2} \quad (2)$$

where  $\hat{C}_C$  and  $\hat{C}_P$  are original signal cepstrum and enhanced signal cepstrum vectors, respectively. In literature, LPC-based objective measures have been utilized, such as log-likelihood ratio (LLR) [14]

$$d_{LLR}(\hat{a}_P, \hat{a}_C) = \log \left( \frac{\hat{a}_P R_C \hat{a}_P^T}{\hat{a}_C R_C \hat{a}_C^T} \right) \quad (3)$$

Where  $\hat{a}_C$  and  $\hat{a}_P$  are LPC vectors of the original and enhanced signals, respectively.  $R_C$  is autocorrelation of original signal.

In [3] composite evaluation is proposed, which was obtained as a correlation between objective and subjective evaluation, by using two merits: correlation coefficient and standard deviation.

Here, a new evaluation measure is proposed by Continuous Wavelet Transform (CWT). This measure is obtained by calculating the differences between CWT of the original signal and the enhanced signal over three levels: low, medium and high. And then, average of standard deviations is obtained

$$d_{CWT} = \frac{\sum_j \sqrt{E[(C_j - \bar{C})^2]}}{3} \quad (4)$$

for  $j = 5, 10$  and  $15$

Where  $C_j = CWT_j(s) - CWT_j(\tilde{s})$  and  $\bar{C}$  is a mean value. The level determination as 5, 10 and 15 is according to the sampling frequency of the speech signal [16]-[17]. These levels present low, medium and high pass bands of the signal frequency. Thus, the utilizing this measure helps studying the difference between filtered and clean signals via three bands, instead of whole signal overlapped bands.

### 3. APPLIED SPEECH ENHANCEMENT METHODS

In this paper we utilize four known speech enhancement method for studying the performance of WNNEM:

#### 1. Discrete Wavelet Filtration Method (DWFM)

This method involves multistage wavelet filtration based on convolution with Reverse

Biorthogonal Wavelets [18]. This method is based on filtration the low frequency and high frequency parts separately, without thresholding (cutting) the values, which leads to lose the essential speech information.

#### 2. Donoho Thresholding Method (DTM)

Donoho and Johnstone in [19] presented soft thresholding function as follows

$$T_S(\lambda, w_k) = \begin{cases} \text{sgn}(w_k)(|w_k| - \lambda) & \text{if } |w_k| > \lambda \\ 0 & \text{if } |w_k| \leq \lambda \end{cases} \quad (5)$$

Where  $w_k$  is the wavelet coefficient, and  $\lambda$  is the universal threshold for WT

$$\lambda = \sigma \sqrt{2 \log(N)} \quad (6)$$

Where  $\sigma = MAD/0.6745$  is the noise level,  $MAD$  is the absolute of median estimated on first scale, and  $N$  is the length a speech frame (de-noised) signal. For Wavelet Packets Transform, the threshold is calculated by

#### 3. Massart Thresholding Method (DTM)

Birgé and Massart in [20] proposed a level-dependent threshold, which can be explained by the following sequent concepts

- $[C, L]$  is the wavelet structure of the decomposed signal to be enhanced (de-noised), at level  $j = \text{length}(L) - 2$ .
- $\alpha$  and  $M$  are real numbers greater than 1.
- $T$  is a vector of length  $j$ ;  $T(i)$  contains the threshold for level  $i$ .
- $N_{KEEP}$  is a vector of length  $j$ ;  $N_{KEEP}(i)$  contains the number of coefficients to be kept at level  $i$ .

The strategy definition:

- 1) For level  $j + 1$ , everything is kept.
- 2) For level  $i$  from 1 to  $j$ , the  $ni$  largest coefficients are kept with  $ni = M (j + 2 - i)^\alpha$ . Typically  $\alpha = 3$  for de-noising. Recommended values for  $M$  are from  $L(1)$  to  $2 * L(1)$ .

#### 4. Kalman Filter Method (KFM)

The time-varying Kalman filter is a generalization of the steady-state filter for time-varying systems or LTI systems with nonstationary noise covariance. More about This filter can be found in [21].

#### 4. FEED FORWARD BACK PROPAGATION NEURAL NETWORK

Feed forward neural networks have been extensively used for a variety of tasks, such as pattern recognition, function dynamical, approximation, data mining, modeling, and time series forecasting [22], [23]. The training of FNN is mostly undertaken by means of the back propagation-based learning algorithms. A number of special kinds of BP learning algorithms have been planned, such as an on line neural network learning algorithm for dealing with time varying inputs, speedy learning algorithms based on gradient descent of neuron space, or the Levenberg–Marquardt algorithm [22]. In this letter, we will develop a back propagation learning algorithm for Feed forward neural networks for speech signal evaluation.

The back propagation algorithm has emerged as the engine for the design of a special class multilayer perceptrons . A multilayer perceptron has an input layer of source nodes and an output layer of neurons (i.e., computation nodes); these two layers connect the network to the outside world. Additionally to these two layers, the multilayer perceptron typically has one or more layers of hidden neurons. The hidden neurons extract significant features contained in the input data. The training of an multilayer perceptrons is usually accomplished by using a back propagation algorithm that involves two phases

- *Forward Phase.* Through this phase the free parameters of the network are fixed, and the input signal is propagated through the network layer by layer. The forward phase finishes with the calculation of an error signal

$$ei = di - yi \quad (9)$$

where  $di$  is the response and  $yi$  is the actual output produced by the network in response to the input  $xi$ .

- *Backward Phase.* Through this second phase, the error signal  $ei$  is propagated through the network in the backward direction, thus the name of the algorithm. It is during this phase that adjustment are applied to the parameters of the network so as to minimize the error  $ei$  in a statistical form. Back-propagation learning may be implemented in one of two basic behavior, as summarized here:

1. *Sequential mode* (also referred to as the on line mode or stochastic mode): In this mode of back propagation learning, adjustments are completed to the free parameters of the network on an example by instance basis. The chronological mode is best suited for classification.

2. *Batch mode:* In this second manner of back propagation learning, adjustments are made to the free parameters of the network on an epoch by epoch

basis, where each epoch consists of the entire set of training examples. The group mode is best suited for nonlinear regression. The back propagation learning algorithm is uncomplicated to apply and computationally competent in that its complexity is linear in the synaptic weights of the network.

#### 5. METHOD

In this paper, we use FFBNN for enhanced signal evaluation by comparing with original free of noise signal. The input P matrix contains N columns of wavelet coefficients; each column presents 2500 wavelet coefficients:

$$P = \begin{bmatrix} I_{F0} & I_{F0} & \dots & I_{F0N} \\ I_{F1} & I_{F1} & \dots & I_{F1N} \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \dots & \cdot \\ I_{F4} & I_{F4} & \dots & I_{F4N} \end{bmatrix} \quad (7)$$

where  $N$  is the column of 2500 wavelet coefficients. To take matching decision, this matrix is given to a FEBNN to be trained with the following binary target for  $N=9$ .

$$T = \begin{bmatrix} 1 & 0 & 1 & 0 \dots \\ 0 & 1 & 1 & 0 \dots \\ 0 & 0 & 0 & 1 \dots \\ 0 & 0 & 0 & 0 \dots \end{bmatrix} \quad (8)$$

To implement FFBNN, we can use matlab neural network toolbox by function newff, tansig transfer function and trainlm backpropagation training function:

```
net=newff(minmax(P),[5 4],{'tansig','trainlm');
```

This command builds a network of three layers: 5 neurons input layer, 5 neurons hidden layer and 4 neurons output layer. After training with the target by  $[net,tr]= \text{train}(P, T)$ ;

We simulate the network outputs (the weights and the biases) with enhanced signal to be evaluated, by  $T\_result=\text{sim}(net, pt)$ ;

Now  $T\_result$  indicates the net output of enhanced signal according to free of noise signal. Now the quality measure is calculated as the distance between  $T\_result$  and the target

$$\text{Error} = \frac{\sum((T-T\_result).^2)}{\sum(T.^2)};$$

## 6. RESULTS AND DISCUSSION

Tested speech signals were recorded via PC-sound card, with a spectral frequency of 4000 Hz and sampling frequency 16000 Hz, over about 2 sec. time duration. For each speaker, the Arabic expression, which sounds "besmeallahalrahmanalraheem", that means in English "In the Name of God the merciful, the compassionate", was recorded 10 times by each speaker. 4 females and 18 males participated in utterances recording. The recording process was provided in normal university office conditions. The experimental part of this research is introduced by utilizing several objective measures such as  $d_{CWT}$ , modified Cepstrum distance

$$Md_{CEP}(\hat{C}_C, \hat{C}_P) = \log \sqrt{2 \sum_{k=1}^p (C_C(k) - C_P(k))^2} \quad (9)$$

and modified LPC-based log-likelihood ratio  $Md_{LLR}$

$$Md_{LLR} = \left| \text{Re} \left( \log \frac{\sum_n^N a_s(n) R_s}{\sum_n^N a_{\bar{s}}(n) R_{\bar{s}}} \right) \right| \quad (10)$$

Where  $a_s(n)$  and  $a_{\bar{s}}(n)$  are LPC of the original and the enhanced signals, respectively.  $R_s$ ,  $R_{\bar{s}}$  are autocorrelation of original and enhanced signals. The modification is done to make the two measures more suitable for our research. Correlation coefficient and MSE are also used.

Table 1: Objective measures results

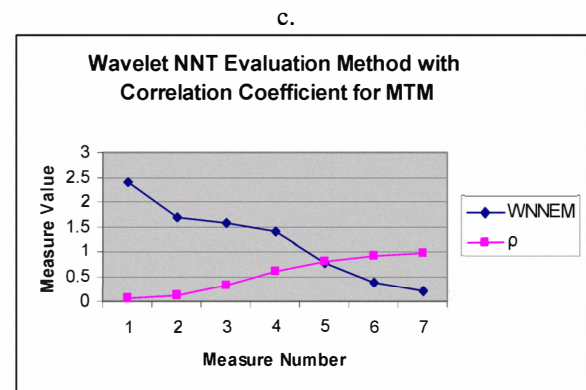
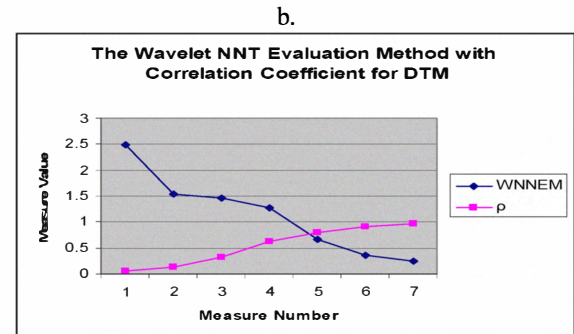
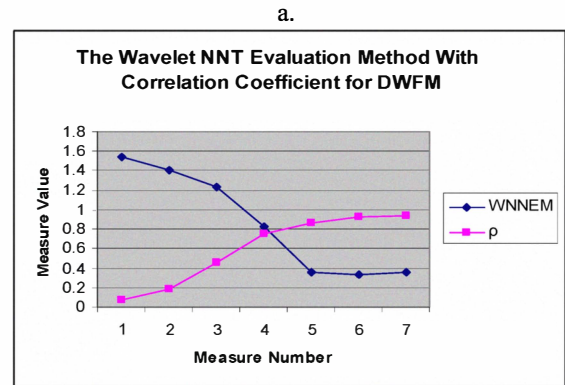
OBJ. EVALUATION	DWFM	DTM	MTM	KFM
<i>SNR</i>	5.1084	2.0276	2.0074	5.5568
<i>P</i>	0.8083	0.7112	0.6991	0.7033
<i>MSE</i>	0.0001	0.0001	0.0002	0.0005
$Md_{CEP}$	0.39	0.47	0.1869	0.3723
$d_{CWT}$	0.0204	0.0212	0.0214	0.0389
<i>WNNEM</i>	0.8465	0.935	0.9458	2.5682

Corrupted Signal SNR=-4.5366 dB

Table 1 contains objective measures results taken for corrupted signal SNR equal to -4.5366 dB. These results were calculated for four enhancement methods mentioned in section three. We can see

clearly the correlation between the conventional objective methods and the proposed method DWFM.

The relation between WNNEM and SNR is presented in table 2. we can see that there is a compatibility between these two measures over four enhancement methods mentioned in section 3. DWFM showed best SNR improvement with best WNNEM (smallest).



d.

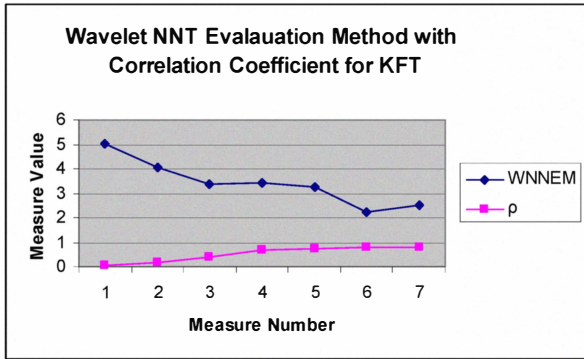


Figure 1: WNNEM with correlation coefficient for enhanced signal by a. DWFM. b. DTM. c. MTM. and KFT

In figure 1 we illustrate the relation between WNNEM and correlation coefficient. These results were calculated for seven SNR levels for corrupted signal, vary from -30 dB to 17 dB, for four enhancement methods mention in section 3. The figures illustrate that there are correct relationship between WNNEM and correlation coefficient, because when correlation coefficient is small then WNNEM as an error is high, but when it is high WNNEM as an error is small.

The main advantage of this method is the capability of error distribution over DWT sub signals. This helps greatly in studying the enhancement methods quality in frequency different sub bands, which is accomplished by using DWT (see table 3). Figure 2 illustrates the WNNEM magnitude for DWT sub signals for DWFM, DTM, MTM and KFT. By this figure we can easily detect the place of deformation by detecting the bigger WNNEM of certain DWT sub signals.

## 7. CONCLUSION

In this paper, Wavelet Neural Network Evaluation method is presented. Feed forward backpropagation neural network is proposed to train the free of noise signal, and then enhanced signal is simulated to the training output results taken for given target. Four known speech enhancement method for studying the performance of Wavelet Neural Network Evaluation method are utilized. The advantage of this method is the evaluation of different band passes of frequency based on wavelet transform by neural network, which is very powerful classification tool. Several objective measures are used to compare the proposed method with. Results proved the validity of the proposed method.

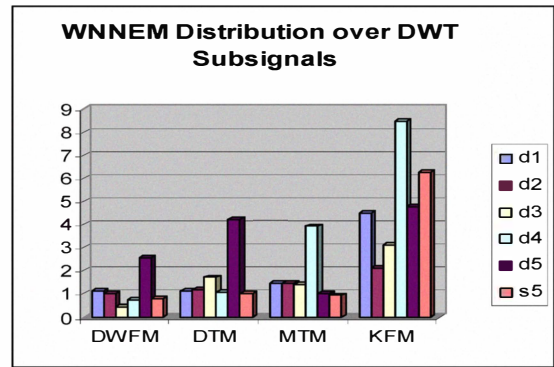


Figure 2: WNNEM over DWT sub signals for DWFM, DTM, MTM and KFT

Table 3: WNNEM for DWT sub signals of DWFM, DTM, MTM and KFT

With Corrupted SNR=-12dB	WNNEM	Improved SNR	d1	WNNEM				
				d2	d3	d4	d5	s5
DWFM	0.7189	0.6166	1.1296	0.9953	0.4493	0.7387	2.5563	0.7684
DTM	1.2545	-1.1405	1.1416	1.1801	1.7046	1.0613	4.2446	1.0128
MTM	1.308	-1.2068	1.4324	1.4668	1.4011	3.9166	1.0258	0.9248
KFM	3.8847	1.2864	4.4956	2.1133	3.1091	8.4692	4.7864	6.2584

## 7. REFERENCES

- [1] Y. Hu and P. Loizou, "Subjective comparison of speech enhancement algorithms," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2006, vol. 1, pp. 153–156.
- [2] Y. Hu and P. Loizou, "Subjective comparison and evaluation of speech enhancement algorithms," *Speech Commun.*, vol. 49, pp. 588–601, 2007.
- [3] Hu Y. and Loizou P., Evaluation of Objective Quality Measures for Speech Enhancement, *IEEE Tran. on Audio, Speech, and Language Processing*, Vol. 16, NO. 1, Jan. (2008).
- [4] S. Quackenbush, T. Barnwell, and M. Clements, *Objective Measures of Speech Quality*. Englewood Cliffs, NJ: Prentice-Hall, 1988.
- [5] "Subjective test methodology for evaluating speech communication systems that include noise suppression algorithm," ITU-T, ITU-T Rec. P. 835, 2003.
- [6] P. Kroon, , W. Kleijn and K. Paliwal, Eds., "Evaluation of speech coders," in *Speech Coding and Synthesis*. New York: Elsevier, 1995, pp. 467–494.
- [7] L. Thorpe and W. Yang, "Performance of current perceptual objective speech quality measures," in *Proc. IEEE Speech Coding Workshop*, 1999, pp. 144–146. [5] T. H. Falk and W. Chan, "Single-ended speech quality measurement using machine learning methods," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 14, no. 6, pp. 1935–1947, Nov. 2006.
- [9] L. Malfait, J. Berger, and M. Kastner, "P.563-the ITU-T standard or single-ended speech quality assessment," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 14, no. 6, pp. 1924–1934, Nov. 2006.
- [10] A. Rix, J. Beerends, M. Hollier, and A. Hekstra, "Perceptual evaluation of speech quality (PESQ)-A new method for speech quality assessment of telephone networks and codecs," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2001, vol. 2, pp. 749–752.
- [11] "Perceptual evaluation of speech quality (PESQ), and objective method for end-to-end speech quality assessment of narrowband telephone networks and speech codecs," ITU, ITU-T Rec. P. 862, 2000.
- [12] R. Kubichek, D. Atkinson, and A. Webster, "Advances in objective voice quality assessment," in *Proc. Global Telecomm. Conf.*, 1991, vol. 3, pp. 1765–1770.
- [13] Turbin and N. Faucheur, "Estimation of speech quality of noise reduced signals," in *Proc. Online Workshop Meas. Speech Audio Quality Netw.*, (2007), [Online]. Available: <http://wireless.feld.cvut.cz/mesaqin2007/contributions.html>.
- [14] Hansen and J. Pellom B, "An effective quality evaluation protocol for speech enhancement algorithms," in *Proc. Int. Conf. Spoken Lang. Process.*, (1998), vol. 7, pp. 2819–2822.
- [15] Klatt D. , "Prediction of perceived phonetic distance from critical band spectra," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, (1982), vol. 7, pp. 1278–1281.
- [16] Daqrouq K. and Abu-Sheikha N.M, Heart Rate Variability Analysis Using Wavelet Transform, *Asian Journal for Information Technology*, Pakistan, Vol. 4, Number4, 2005.
- [17] Daqrouq K. and Abu-Isbeih Ibrahim N., Arrhythmia Detection Using Wavelet transform, *IEEE Region 8, EUROCON 2007*, Warsaw, Poland, Sept., 2007.
- [18] Khaled Daqrouq, Abdel-Rahman Al-Qawasmi, The Study of Wavelet Filters Speech Enhancement Method, *MIC-CCA2009*, 26-28 Oct. 2009
- [19] Donoho, D., Johnstone, I., (1994). Ideal spatial adaptation by wavelet shrinkage. *Biometrika* 81, 425–455.
- [20] Birgé, L.; P. Massart (1997), "From model selection to adaptive estimation," in D. Pollard (ed), *Festschrift for L. Le Cam*, Springer, pp. 55-88.
- [21] Grimble, M.J., *Robust Industrial Control: Optimal Design Approach for Polynomial Systems*, Prentice Hall, 1994, p. 261 and pp. 443-456.
- [22] R. Parisi, E. D. Di Claudio, G. Orlandi, and B. D. Rao, "A generalized learning paradigm exploiting the structure of feedforward neural networks," *IEEE Trans. Neural Networks*, vol. 7, pp. 1450–1459, Nov. 1996.
- [23] M. T. Hagan and M. B. Menhaj, "Training feedforward neural networks with the Marquardt algorithm," *IEEE Trans. Neural Networks*, vol. 5, pp. 989–993, Nov. 1994