

Deep Neural Networks for Shimmer Approximation in Synthesized Audio Signal

Mario Alejandro García ✉ and Eduardo Atilio Destéfánis

Universidad Tecnológica Nacional Facultad Regional Córdoba, Argentina
mgarcia@frc.utn.edu.ar

Abstract. Shimmer is a classical acoustic measure of the amplitude perturbation in a signal. This kind of variation in the human voice allows to characterize some properties, not only of the voice itself, but of the person who speaks. During the last years deep learning techniques have become the state of the art for recognition tasks on the voice. In this work the relationship between shimmer and deep neural networks is analyzed. A deep learning model is created. It is able to approximate shimmer value of a simple synthesized audio signal (stationary and without formants) taking the spectrogram as input feature. It is concluded firstly, that for this kind of synthesized signal, a neural network like the one we proposed can approximate shimmer, and secondly, that the convolution layers can be designed in order to preserve the information of shimmer and transmit it to the following layers.

Keywords: shimmer, voice quality, deep learning, deep neural network, convolutional neural network

1 Introduction

Shimmer is a classical acoustic measure of the amplitude perturbation of a signal. This kind of variations in the human voice allows to characterize some properties, not only of the voice itself, but on the person who speaks [1].

Shimmer value is associated to voice quality [2–7], state of mind [8–13], age [14] and gender [15] of people. There are many research works that use shimmer (among other measures) with goals ranging from pathologies detection [6, 16, 17] to the improvement of human-machine interfaces through the estimation of the intensionality of a spoken phrase [18]. Regarding synthesized voices, Yamasaki et al. show in [19] that a certain shimmer level increases the degree of naturalness.

The application of deep learning techniques is the state of the art in automated audio analysis, with the detection of pronounced phonemes and the identification of the person that speaks as main objectives [20–26], but also used to detect emotions, age, gender, etc. [27–33].

Classifiers based on neural networks can be divided into two groups according to the type of input features, those using previously calculated acoustic measures [10, 14] and those using raw audio [22, 24, 25] or spectral data [21–23, 28–31, 34, 35]. In [26] a hybrid approach is applied by adding shimmer and other measures

to improve the recognition achieved with spectral data. It is important to clarify, for first group of classifiers, that shimmer calculation has a major complication, it depends on the previous detection of the fundamental frequency (f_0) of vocal cords vibration. It is difficult to estimate f_0 in pathological voices [36, 37]. The estimation of the actual f_0 value is still a research topic [36–40]. Regarding the second group of classifiers, it is not possible to know whether the outputs are influenced by the shimmer value of the signal.

1.1 Objectives

The objective of this work is to make an estimation of the shimmer value in a synthesized audio signal through a neural model. The neural network must combine convolutional layers and feed forward layers. The inputs of the neural model will be the spectral values of the signal.

Main contributions of developing a neural network that estimates shimmer from spectral features of an audio signal are, on the one hand, the procurement of a f_0 independent shimmer calculation method, and on the other, to answer the question about the extent to which amplitude disturbances of the original audio can influence the output of a deep learning model with raw audio or spectral data input. In other words, how much shimmer information is preserved to the last layers of the model.

1.2 Shimmer Calculation

There are different versions of shimmer. The most important difference between them is the window size (number of f_0 cycles) used for the calculation. Some versions can be seen in [41].

Chosen version of shimmer for this work is the proposed by Klingholz and Martin [42], also known as *Relative Shimmer*.

Relative Shimmer, hereinafter referred to as "shimmer", is a way of measuring cycle-to-cycle amplitude perturbations of the fundamental frequency of a signal. It is shown as a *perturbation/total amplitude* relation.

$$\text{shimmer} = \frac{\frac{1}{N-1} \sum_{i=1}^{N-1} |A_i - A_{i+1}|}{\frac{1}{N} \sum_{i=1}^N |A_i|} \quad (1)$$

where N is the number of periods of f_0 in the signal and A_i is the maximum amplitude into i period.

2 Methods and Materials

2.1 Neural Models

Deep learning models with ascending complexity were generated for problems of shimmer approximation. First, shimmer was approximated for f_0 variable, k

constant and f_{mod} constant. Then, shimmer was approximated for f_0 variable, k variable, and f_{mod} constant. Finally, a model was found to approximate shimmer with f_0 , k and f_{mod} variable.

In all cases, spectral audio data (instead of raw audio) were used as input features. There are two reasons, the improvement of training performance due the dimension reduction, and the similarity with human auditory system, where spectral decomposition is performed in the basilar membrane of the cochlea and not by the neurons in auditory cortex [43].

2.2 Data

Audio. Audio data without harmonics was generated. As in [1] the amplitude modulation of human voice was approximated by a sinusoidal wave. The expression of each audio signal $y(t)$ was:

$$y(t) = \frac{1}{1+k} \sin(\alpha + 2t\pi f_0)(1 + k \sin(\beta + 2t\pi f_{\text{mod}})) \quad (2)$$

where t is time [sec], f_0 is the frequency of vocal fold vibration [Hz], f_{mod} is the modulatory frequency [Hz], k is the constant of the amplitude modulator sensibility, α and β are constant to handle the phase of the signal to be modulated and the modulating signal respectively.

For training, test and validation data generation, random values were taken with uniform distribution. f_0 got values in [200, 1000] Hz range, f_{mod} in [5, 10] Hz, k in [0, 0.4], α and β in $[0, 2\pi]$.

250 ms of audio generated with $f_0 = 200$ Hz, $f_{\text{mod}} = 8$ Hz and $k = 0.4$. are shown in Fig. 1.

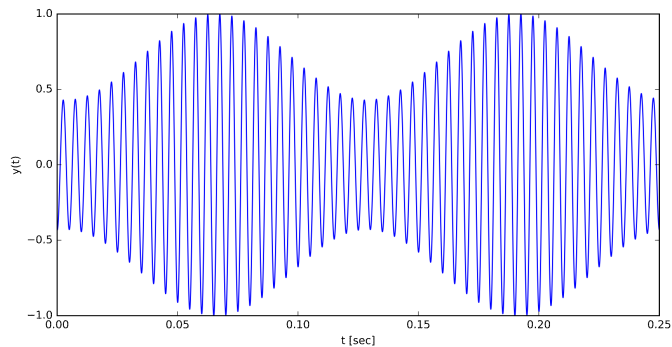


Fig. 1. Generated audio for $f_0 = 200$ Hz, $f_{\text{mod}} = 8$ Hz and $k = 0.4$

Training data. Three datasets were created to train each model, a training dataset with 2500 elements, a testing dataset with 500 elements and a validation dataset with 500 elements. Each element is composed by shimmer (Eq.(1)) value to be estimated and the spectrogram of generated audio.

Due to the fact that f_0 is known at the time of audio generation, shimmer value can be accurately calculated.

The spectrogram is calculated on 2 seconds of 44100 samples/sec audio. A *Tukey*(0.25) window of width = 256 was used, which determines a 129 x 393 (frequency/time) shape structure that contains the signal spectral density.

Fig. 2 shows values of the second, third and fourth rows (index 1 to 3) of the spectrogram of signal in Fig. 1.

Spectrograms data and shimmer data were scaled into the range [0,1].

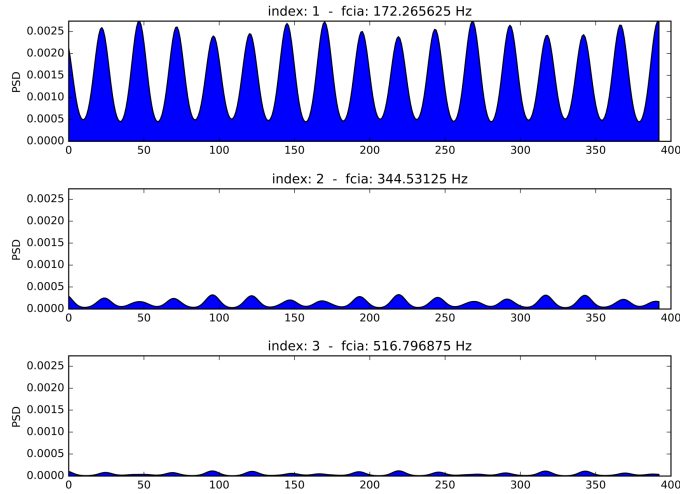


Fig. 2. Three rows with higher average value of Power Spectral Density (PSD) in spectrogram of audio generated with $f_0 = 200$ Hz, $f_{\text{mod}} = 8$ Hz and $k = 0.4$

3 Results

An initial analysis was performed with f_0 , f_{mod} and k known data. It was found that a neural network with dense connections is able to calculate shimmer value with high precision if it gets f_0 , f_{mod} and k as input features. Optimal structure of this network was empirically found. This is a three layer network, two layers of 20 neurons with $\tanh()$ activation function and a linear neuron as output. In next models, convolution layers are used at the initial part of the network, and then, dense layers with 20, 20 and 1 neurons. The function of convolution layers is to calculate f_0 , f_{mod} and k values in order to dense layers calculate shimmer.

It was noted that only the first 15 rows of the spectrogram (lower frequencies) would have significant information. Then, only for the scope proposed in this work, the rest of the frequencies were deleted. Spectrogram shape changes from 129×393 to 15×393 . This provides a important performance improvement.

3.1 Shimmer Approximation with f_0 Variable

Without harmonics, f_0 calculation from spectrogram is easy, it is enough to obtain the energy average value weighted by the frequency that each spectrogram row represents. As expected, a network such as Fig. 3, where each complete row of the spectrogram is connected to one neuron of an *Average pooling* layer, is able to perform the weighted average of frequencies and calculate the shimmer value in densely connected layers. Tests were performed with f_0 in $[200, 1000]$ Hz range, $k = 0.4$ and $f_{\text{mod}} = 8$ Hz. A satisfactory approximation was achieved, with a mean square error (MSE) $< 10^{-4}$.

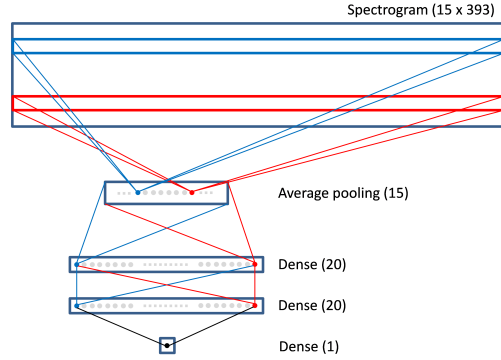


Fig. 3. Shimmer approximation model on signals with f_0 in $[200, 1000]$ Hz range, $f_{\text{mod}} = 8$ Hz and $k = 0.4$. Each neuron in the *Average pooling* layer has a complete frequency of the spectrogram as its visual field. The activation function of hidden dense layers neurons is $\tanh()$ and the output neuron is linear.

3.2 Shimmer Approximation with f_0 and k Variable

The value of k inversely affects the area under the energy curve of the spectrogram. Therefore, information about k value can be obtained through the energy average of the spectrogram. The model shown in previous section preserves the necessary information to estimate the energy average. Tests were performed with audio data for f_0 in $[200, 1000]$ Hz range, k in $[0, 0.4]$ range and $f_{\text{mod}} = 8$ Hz. Results were satisfactory again. The model approximates shimmer with an MSE $< 10^{-4}$.

3.3 Shimmer Approximation with f_0 , k and f_{mod} Variable

For f_0 in the range [200, 1000] Hz, k in the range [0, 0.4] and f_{mod} in the range [5, 10] Hz it was necessary to create a more complex model than the previous one. Shimmer depends on the modulation frequency, so a new transformation is necessary. The first one was the transformation from time domain to frequency domain (spectrogram). The new (second) transformation is performed in a convolution layer at the initial part of the model (Fig. 4).

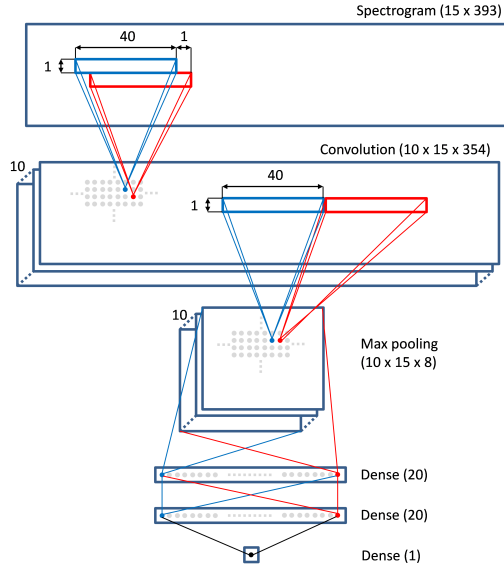


Fig. 4. Shimmer approximation model on signals with f_0 in the range [200, 1000] Hz, k in the range [0, 0.4] and f_{mod} in the range [5, 10] Hz. The shape of convolutional layer windows is 1 x 40, strides 1 x 1. Convolutional layer has 10 sub-layers. The shape of *max pooling* layer windows is 1 x 40, strides 1 x 40. The network finishes with three dense layers of 20, 20 and 1 neurons.

Convolutional layer. Each convolution layer neuron is connected to spectrogram through a $height = 1$ and $width = 40$ window. Convolution is performed on a single frequency ($height = 1$) in order to the f_0 detail level needed in the dense layers is not lost. 40-element width is the minimum required to hold a cycle of $\min(f_{\text{mod}})$. The number of elements of the spectrogram per modulation cycle (C) for a spectrogram of width W_s and an audio length L is:

$$C = \frac{W_s}{L \times \min(f_{\text{mod}})} = \frac{393 \text{ elements}}{2 \text{ sec} \times 5 \text{ Hz}} = 39.3 \text{ elements/cycle}$$

The window displacement in both directions is 1 step. This implies that on the frequency dimension there is no overlap, and in time dimension there are 39 overlapping elements between the windows of adjacent neurons. Finally, according to these definitions, the shape of each convolution filter or sub-layer is 15×354 . The convolution layer consists of 10 sub-layers. This amount is a compromise between performance and the detail level of f_{mod} on the information sent to dense layers. Neurons of this layer have linear activation function. Weights are initialized with orthogonal random values. An attempt was made to initialize them with wavelet families for sinusoidal waves between 5Hz y 10Hz, but no improvement was achieved on the prediction accuracy.

Max pooling layer. The neurons in the *max pooling* layer have a 1×40 window size on the convolutional layer. Again, *height* = 1 allows f_0 information be able to be transmitted to dense layers with no losing details. The 40-element width extends the visual field of this layer neurons to 2 cycles of $\min(f_{\text{mod}})$ on the spectrogram. In this way, the output value is invariant to the modulation signal translations. There is no overlap between the windows, so the size of each of the 10 sub-layers is 15×8 neurons.

The outputs of *max pooling* layer are connected to three layers with dense connections equal to those of the previous model.

For this model, 20 training tests were performed. The size of training dataset was 2500 elements. In all cases, results were compared with a test dataset (500 elements) during training and a validation dataset (500 elements) at the end. The best result, with 150 training cycles, obtained a $\text{MSE} = 5.8 \times 10^{-5}$ on the test dataset. In Fig. 5 expected and calculated shimmer values are displayed in ascending order for the 500 elements of test dataset.

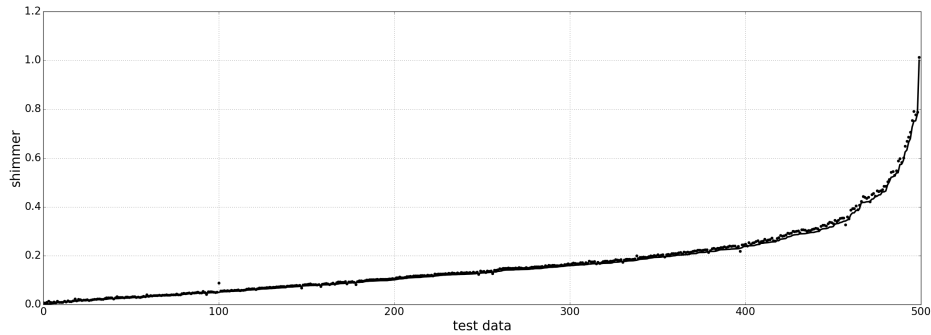


Fig. 5. Normalized shimmer. Expected (line) vs. calculated (dots) for elements in testing dataset (in ascending order of shimmer value).

4 Conclusion

It was verified that, for simple audio signal modulated in amplitude by a sinusoidal wave, with variable parameters of fundamental frequency, modulating frequency and modulation sensitivity, it is possible to obtain a neural model able to approximate the value of shimmer.

Under the conditions established in this work, it is possible to calculate shimmer without knowing f_0 . Moreover, it can be affirmed that if the first layers of a deep neural network respects the structure of the second presented model, this neural network is able to use the value of shimmer, internally calculated, to perform classifications and other approximations.

5 Future Works

It is planned to extend the analysis, first by expanding the ranges of f_0 y f_{mod} , then adding harmonics and noise to the synthesized signals. Finally, it is planned to analyze the behavior of deep learning models for shimmer calculation on natural voices.

References

1. Jafari, M., Till, J.A., Law-Till, C.B.: Interactive effects of local smoothing window size and fundamental frequency on shimmer calculation. *Journal of Voice* **7**(3) (1993) 235–241
2. Nieto, R.G., Marín-Hurtado, J.I., Capacho-Valbuena, L.M., Suarez, A.A., Bolaños, E.A.B.: Pattern recognition of hypernasality in voice of patients with cleft and lip palate. In: *Image, Signal Processing and Artificial Vision (STSIVA), 2014 XIX Symposium on, IEEE* (2014) 1–5
3. Holi, M.S., et al.: A hybrid model for neurological disordered voice classification using time and frequency domain features. *Artificial Intelligence Research* **5**(1) (2015) 87
4. Freitas, S.V., Pestana, P.M., Almeida, V., Ferreira, A.: Integrating voice evaluation: correlation between acoustic and audio-perceptual measures. *Journal of Voice* **29**(3) (2015) 390–e1
5. Little, M.A., Costello, D.A., Harries, M.L.: Objective dysphonia quantification in vocal fold paralysis: comparing nonlinear with classical measures. *Journal of Voice* **25**(1) (2011) 21–31
6. Lopes, L.W., Simões, L.B., da Silva, J.D., da Silva Evangelista, D., e Ugulino, A.C.d.N., Silva, P.O.C., Vieira, V.J.D.: Accuracy of acoustic analysis measurements in the evaluation of patients with different laryngeal diagnoses. *Journal of Voice* **31**(3) (2017) 382–e15
7. Hillenbrand, J.: Perception of aperiodicities in synthetically generated voices. *The Journal of the Acoustical Society of America* **83**(6) (1988) 2361–2371
8. Li, X., Tao, J., Johnson, M.T., Soltis, J., Savage, A., Leong, K.M., Newman, J.D.: Stress and emotion classification using jitter and shimmer features. In: *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on. Volume 4., IEEE* (2007) IV–1081

9. Kotti, M., Stylianou, Y.: Effective emotion recognition in movie audio tracks. In: *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*, IEEE (2017) 5120–5124
10. Jacob, A.: Speech emotion recognition based on minimal voice quality features. In: *Communication and Signal Processing (ICCSP), 2016 International Conference on*, IEEE (2016) 0886–0890
11. Sondhi, S., Vijay, R., Khan, M., Salhan, A.K.: Voice analysis for detection of deception. In: *Knowledge, Information and Creativity Support Systems (KICSS), 2016 11th International Conference on*, IEEE (2016) 1–6
12. Palo, H.K., Mohanty, M.N., Chandra, M.: Sad state analysis of speech signals using different clustering algorithm. In: *Next Generation Computing Technologies (NGCT), 2016 2nd International Conference on*, IEEE (2016) 714–718
13. Schuller, B., Steidl, S., Batliner, A.: The interspeech 2009 emotion challenge. In: *Tenth Annual Conference of the International Speech Communication Association*. (2009)
14. Kim, H.J., Bae, K., Yoon, H.S.: Age and gender classification for a home-robot service. In: *Robot and Human interactive Communication, 2007. RO-MAN 2007. The 16th IEEE International Symposium on*, IEEE (2007) 122–126
15. Teixeira, J.P., Fernandes, P.O.: Jitter, shimmer and hnr classification within gender, tones and vowels in healthy voices. *Procedia Technology* **16** (2014) 1228–1237
16. Tsanas, A., Little, M.A., Fox, C., Ramig, L.O.: Objective automatic assessment of rehabilitative speech treatment in parkinson’s disease. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* **22**(1) (2014) 181–190
17. Gómez-Coello, A., Valadez-Jiménez, V.M., Cisneros, B., Carrillo-Mora, P., Parra-Cárdenas, M., Hernández-Hernández, O., Magaña, J.J.: Voice alterations in patients with spinocerebellar ataxia type 7 (sca7): Clinical-genetic correlations. *Journal of Voice* **31**(1) (2017) 123–e1
18. Kotti, M., Paternò, F.: Speaker-independent emotion recognition exploiting a psychologically-inspired binary cascade classification schema. *International journal of speech technology* **15**(2) (2012) 131–150
19. Yamasaki, R., Montagnoli, A., Murano, E.Z., Gebrim, E., Hachiya, A., da Silva, J.V.L., Behlau, M., Tsuji, D.: Perturbation measurements on the degree of naturalness of synthesized vowels. *Journal of Voice* **31**(3) (2017) 389–e1
20. Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A.r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T.N., et al.: Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine* **29**(6) (2012) 82–97
21. Mitra, V., Sivaraman, G., Nam, H., Espy-Wilson, C., Saltzman, E., Tiede, M.: Hybrid convolutional neural networks for articulatory and acoustic information based speech recognition. *Speech Communication* **89** (2017) 103–112
22. Collobert, R., Puhersch, C., Synnaeve, G.: Wav2letter: an end-to-end convnet-based speech recognition system. *arXiv preprint arXiv:1609.03193* (2016)
23. Amodei, D., Ananthanarayanan, S., Anubhai, R., Bai, J., Battenberg, E., Case, C., Casper, J., Catanzaro, B., Cheng, Q., Chen, G., et al.: Deep speech 2: End-to-end speech recognition in english and mandarin. In: *International Conference on Machine Learning*. (2016) 173–182
24. Palaz, D., Collobert, R., et al.: Analysis of cnn-based speech recognition system using raw speech as input. *Technical report, Idiap* (2015)
25. Sainath, T.N., Kingsbury, B., Mohamed, A.r., Ramabhadran, B.: Learning filter banks within a deep neural network framework. In: *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*, IEEE (2013) 297–302

26. Farrús, M.: Jitter and shimmer measurements for speaker recognition. In: 8th Annual Conference of the International Speech Communication Association; 2007 Aug. 27-31; Antwerp (Belgium).[place unknown]: ISCA; 2007. p. 778-81., International Speech Communication Association (ISCA) (2007)
27. Gu, Y., Li, X., Chen, S., Zhang, J., Marsic, I.: Speech intention classification with multimodal deep learning. In: Canadian Conference on Artificial Intelligence, Springer (2017) 260–271
28. Chang, J., Scherer, S.: Learning representations of emotional speech with deep convolutional generative adversarial networks. arXiv preprint arXiv:1705.02394 (2017)
29. Ghosh, S., Laksana, E., Morency, L.P., Scherer, S.: Representation learning for speech emotion recognition. In: INTERSPEECH. (2016) 3603–3607
30. Mao, Q., Dong, M., Huang, Z., Zhan, Y.: Learning salient features for speech emotion recognition using convolutional neural networks. *IEEE Transactions on Multimedia* **16**(8) (2014) 2203–2213
31. Ma, X., Yang, H., Chen, Q., Huang, D., Wang, Y.: Depaudionet: An efficient deep model for audio based depression classification. In: Proceedings of the 6th International Workshop on Audio/Visual Emotion Challenge, ACM (2016) 35–42
32. Abumalouh, A., Qawaqneh, Z., Barkana, B.D.: Deep neural network combined posteriors for speakers' age and gender classification. In: Industrial Electronics, Technology & Automation (CT-IETA), Annual Connecticut Conference on, IEEE (2016) 1–5
33. Qawaqneh, Z., Mallouh, A.A., Barkana, B.D.: Deep neural network framework and transformed mfccs for speaker's age and gender classification. *Knowledge-Based Systems* **115** (2017) 5–14
34. Liu, Y., Wang, X., Hang, Y., He, L., Yin, H., Liu, C.: Hypemasality detection in cleft palate speech based on natural computation. In: Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD), 2016 12th International Conference on, IEEE (2016) 523–528
35. Cummins, N., Epps, J., Ambikairajah, E.: Spectro-temporal analysis of speech affected by depression and psychomotor retardation. In: Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on, IEEE (2013) 7542–7546
36. Teixeira, J.P., Gonçalves, A.: Algorithm for jitter and shimmer measurement in pathologic voices. *Procedia Computer Science* **100** (2016) 271–279
37. Shahnaz, C., Zhu, W.P., Ahmad, M.O.: A new technique for the estimation of jitter and shimmer of voiced speech signal. In: Electrical and Computer Engineering, 2006. CCECE'06. Canadian Conference on, IEEE (2006) 2112–2115
38. Dong, B.: Characterizing resonant component in speech: A different view of tracking fundamental frequency. *Mechanical Systems and Signal Processing* **88** (2017) 318–333
39. Liu, B., Tao, J., Zhang, D., Zheng, Y.: A novel pitch extraction based on jointly trained deep blstm recurrent neural networks with bottleneck features. In: Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on, IEEE (2017) 336–340
40. Schlotthauer, G., Torres, M.E., Rufiner, H.L.: A new algorithm for instantaneous f₀ speech extraction based on ensemble empirical mode decomposition. In: Signal Processing Conference, 2009 17th European, IEEE (2009) 2347–2351
41. Buder, E.H.: Acoustic analysis of voice quality: A tabulation of algorithms 1902–1990. *Voice quality measurement* (2000) 119–244

42. Klingholz, F., Martin, F.: Quantitative spectral evaluation of shimmer and jitter. *J Speech Hear Res* **28**(2) (1985) 169–174
43. Schnupp, J., Nelken, I., King, A.: *Auditory neuroscience: Making sense of sound*. MIT press (2011)