# Chapter 8

# Classifying Recurrent Dynamics on Emotional Speech Signals

**Sudhangshu Sarkar*** and **Anilesh Dey†**
*Department of Electrical Engineering, Narula Institute of Technology, Kolkata, India,
†Department of Electronics and Communication Engineering, Narula Institute of Technology, Kolkata, India

## 8.1 Introduction

Speech signal processing is a vast field of study, which contributes highly to human computer interaction (HCI) studies. Speech is simply the best carrier of information in human communication systems. It contains a lot of information apart from the verbal message, that is, speaker identification, speaker's age, sex, locality, emotion, etc. Recognition of emotion contained in a speech signal is one of the fastest growing areas of interest in HCI study. Emotion in a speech plays an important role in expressing feelings. Based on different emotions, human speak in different ways, and the characteristics of speech changes, accordingly. Humans don't need practice to recognize the emotional state of a speaker; it comes naturally [1]. However, it is a complex process when it's implemented in a machine.

In this regard, a literary review of the past studies on speech-based emotion recognition systems was made. Researchers have proposed many techniques for emotion-based speech recognition [2, 3]. Schuller et al. [4] reported the continuous use of hidden Markov models (HMM) for speech-emotion recognition. The same group further extended their work in 2004 by combining acoustic features with linguistic data for a healthy emotion detection using support vector machine (SVM) [5]. Lin et al. [6] used the aforementioned two methods, namely HMM and SVM for classification of five dissimilar states of emotion, that is, annoyance, pleasure, sorrow, shock, and impartial emotion. Lalitha et al. [7] reported the use of time-domain speech features like pitch and prosody for recognition of seven different emotional states. Kamal et al. [8] predicted protein structures from images using HMM. Chapman Kolmogrov. Dey et al. [9,10] analyzed the progressiveness of acoustic waves in biomedical technology.

**139**

Identifying suitable features that characterize different emotions is an important process for developing a speech emotion recognition (SER) system [11]. Altrov et al. [12] identified the power of verbal communication and civilization on the accepting of language emotions. An effective discussion mainly depends on how we communicate our own emotions, how we recognize those of others, and how sufficient our reaction is to their emotions. Cowie et al. [13,14] illustrated that emotions have an essential role in our lives since they are typically present in everyday communication. Park and Sim [15] showed emotion detection by DRNN. Their paper found that pitch was a significant component in the identification of emotion. Therefore the basic accurate detection acoustical features [16,17] were analyzed for speech sound with emotion. The study of event-related potentials (ERPs) is recognized as a useful technique for exploring intelligent mechanisms of processing emotional speech [18]. Tao et al. [19] attempts to create exciting dialogue via "strong," "average," and "weak" classifications using various models like a linear modification model (LMM), a Gaussian mixture model (GMM), a classification and regression tree model (CART). Kang and Li [19] analyzed neutral-emotional speech by using prosody conversion. In 2010, Wu et al. [20] presented an advance to hierarchical prosody translation for an exciting language mixture. Jia et al. [21] adopted "Emotional Audio Visual Speech Synthesis Based on PAD," while Dey and Ashour [22, 23] discussed arrival estimation of localized multispeech sources. An emotional text-to-speech system [24] is required for emotion-based speech recognition. Neural networks [25–27] have exhibited remarkable success to link the responsive space in communication signals. The direction of speech resources on a localized level has been eminently described by Dey and Ashour [28–30]. least squares regression [31] is one of the noted methods for speech emotion recognition. An ideal scientific SER system would be one that can develop real life and loud talking to recognize different state of emotions. In this paper, we have attempted to classify the recurrent dynamics of two different emotions, namely anger and normal, with the help of recurrence plot, phase space plot, and recurrence based parameters. The investigation was performed in both noise-free and on noisy environment to establish the suitability of the proposed method.

## 8.2  Data Collection and Processing

A healthy male volunteer (age 23 years old) was asked to participate in the study. He was informed about the details of the study, and a written consent to participate was obtained. Two types of speech signals were acquired from the volunteer using microphone Behringer C-1U, when he uttered eight different sentences (in the Bengali language) in angry and normal emotion. In order to process the speech signals, the sampling frequency was taken as 16 KHz in using Audacity version 1.3.6. in the Electronics and Communication Engineering Department, Narula Institute of Technology, Agarpara, Kolkata.

## 8.3   Research Methodology

Phase space approach was used to investigate the nonlinear properties of the speech signals. A phase space was reconstructed for each speech signal with appropriate time delay and proper embedding measurement.

### 8.3.1   Phase Space Reconstruction

The condition of a dynamical method is able to be illustrated in a space called phase space. A phase space is a multidimensional space, in which every point correlates with one state of the dynamical system [32]. The path traced by the phase space diagram of a system over time describes its evolution from an initial state. This is known as phase space trajectory.

The basic problem is that the information about all the variables governing the system is usually not obtained from its time series. Most of the time the series is single valued. Although numerous concurrent measurements can be performed; they may not coat each degree of freedom of the arrangement. Though the use of the time-delay embedding theorem [33] allows the recovery of lost information, and it becomes possible to construct the phase space diagram of a scheme from its period sequence. This method of phase space reconstruction requires the determination of the principles of suitable time delay $\tau$ and proper embedding measurement m.

Determining the most favorable value of time delay $\tau$ for the phase space reformation, for any time series $\{x(t)\}_{t=1}^{N}$ at a given state $x(t)$, $\tau$ is one of the proper values of time delay, which divulges utmost novel information through dimension at $x(t+\tau)$. The auto mutual information (AMI) technique [34] is usually adapted the proper value of $\tau$. The AMI of a time series for a given $\tau$ is calculated using Eq. (8.1) [35]. The optimal value of $\tau$ is that one for which AMI($\tau$) reaches its first minimum [34].

$$\text{AMI}(\tau) = \sum_{t=1}^{N-\tau} P[x(t), x(t+\tau)] \log \left( \frac{P[x(t), x(t+\tau)]}{P[x(t)]P[x(t+\tau)]} \right) \qquad (8.1)$$

where $\tau = [1, 2, ..., N-1]$ and $P[\ ]$ denotes the probability.

Embedding dimension is a measure of the least element of the phase space of the reconstructed characteristic of a dynamical system [36, 37]. Kennel et al. [38] have anticipated the method of false nearest neighbor (FNN) to determine the minimum satisfactory embedding dimension $m$. The FNN algorithm can be described as follows.

For every point $\vec{R}_i$ in the time series, its adjacent neighbor $\vec{R}_j$ is searched in an $m$- dimensional space. The space $\left\| \vec{R}_i - \vec{R}_j \right\|$ is calculated. Both the points are iterated and $\vec{R}_i$ is computed as given in Eq. (8.2).

$$R_i = \frac{|R_{i+1} - R_{j+1}|}{\left\| \vec{R_i} - \vec{R_j} \right\|} \tag{8.2}$$

If the computed value of $\vec{R_i}$ goes beyond a specified heuristic verge $\vec{R_t}$, this point is regarded as having a fake nearby neighbor. The minimal embedding length is obtained when the percentage of FNN at a given dimension reaches zero.

### 8.3.2   Recurrence Plot Analysis

Recurrence plot (RP) is one of the efficient graphical methods designed to find the hidden nonlinear structure of the phase spaces, introduced by Eckmann et al. [39]. For any two points $x_i$, $x_j$ in a phase space, the distance among $x_i$ and $x_j$ can be calculated by $\|x_i - x_j\|$. Then, the recurrence between two points $x_i$ and $x_j$ is given by.

$$R_{ij} = \Theta\left(\varepsilon - \|x_i - x_j\|\right) \tag{8.3}$$

where $\Theta$ represents Heaviside function.

From the definition, it follows that the entries in the matrix $(R_{ij})_{N \times N}$ ($N$ being the span of the trajectory of the phase space) are either 1 or 0. The number "1" is represented by a black dot. On the other hand, "0" is represented by a white dot. So, an RP is a visual representation of a phase space by two colors. From the structure of the RP, various dynamical patterns of a complex dynamic can be described, such as periodicity, quasiperiodicity, noise effect, nonstationary behavior, and a chaotic nature. It indicates that classification between two different dynamics can be made by RP analysis. Fig. 8.1(A) and (B) show the recurrence plot of a speech signal in angry emotion and normal emotion, respectively.

If two points $x_i$ and $x_j$ are recurrent, we say that there is an isometry. Two points $x_i$ and $x_j$ in a phase space are said to be in consecutive isometry if.

$$R_{ij} = \Theta\left(\varepsilon - \|x_{i+L} - x_{j+L}\|\right) \tag{8.4}$$

where $L \in Z^+$.

Since periodicity and aperiodicity of a phase space are the reflection of isometry, so the nature of the dynamics can be described by it. In fact, complexity of the RP decreases as the consecutive isometry increases [40, 41].

## 8.4   Numerical Experiments and Results

In this work, speech signals in angry and normal emotion were analyzed using RP analysis for two cases, that is, noise free and noisy conditions. Two informative parameters, namely isometry and consecutive isometry were calculated for each case. For both cases, proper time delay and embedding dimension were
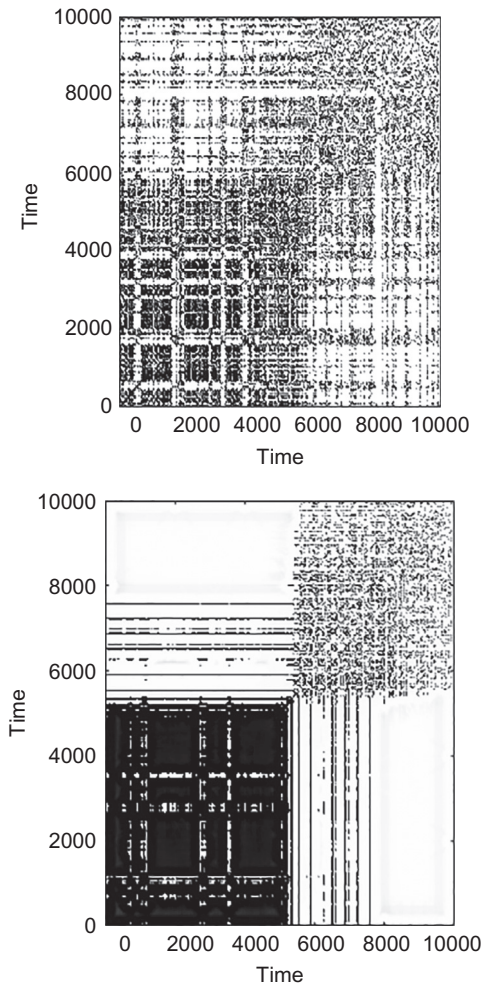
**FIG. 8.1** (A) Recurrence plot of angry emotional speech signal with proper time delay and embedding dimension. (B) Recurrence plot of normal emotional speech signal with proper time delay and embedding dimension.

recalculated. Fig. 8.2(A) shows the time delay for a speech signal in both the angry and normal emotion, which are represented by red and blue lines respectively. In order to calculate the probability, 17 bins were considered.

From the figure, it is observed that the AMI was minimum at $\tau = 60$ in the case of normal speech, whereas it was minimum at $\tau = 43$ for the angry speech signal. It suggested that the optimal time delay was different for the same speech in two different emotions. In fact, it was higher in the case of normal speech than that of the speech in angry emotion.
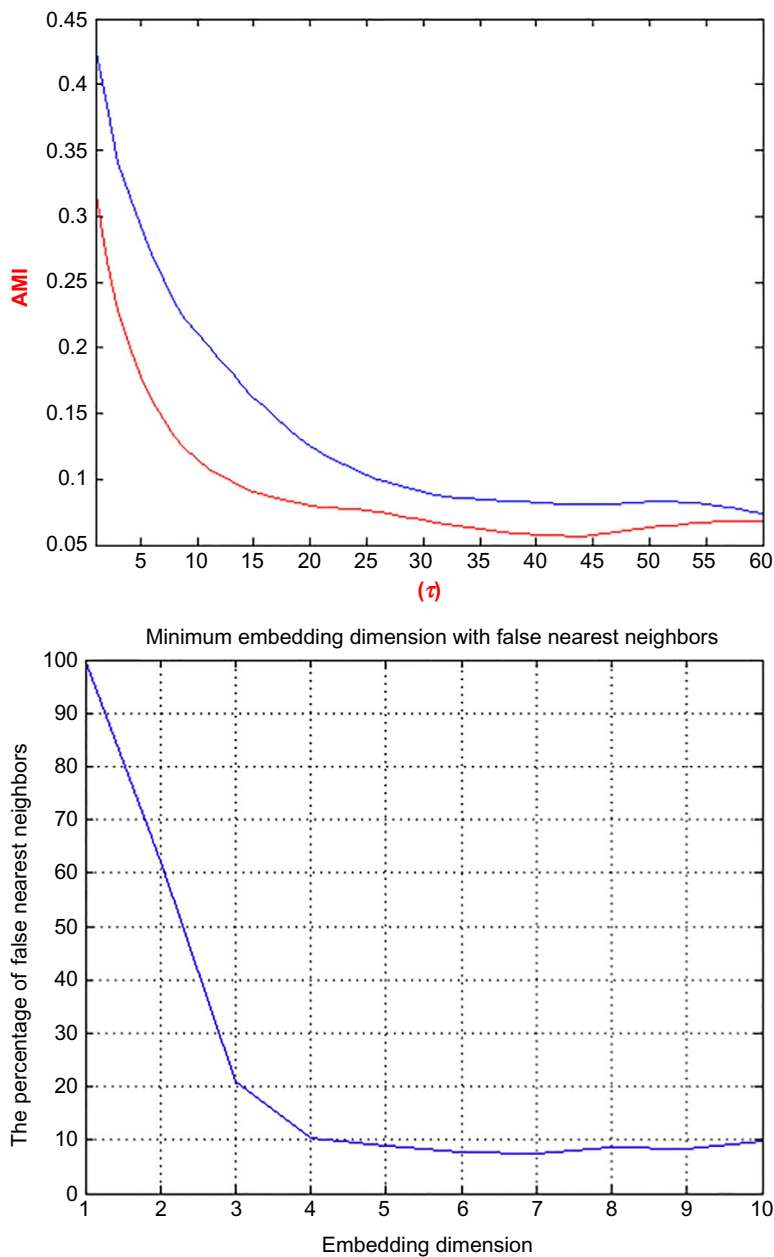
**FIG. 8.2**  (A) Fluctuation of AMI with time delay $\tau \in [1, 60]$. (B) Fluctuation of FNN for a speech signal in the angry emotion with embedding. (C) A phase space diagram for speech signals in two different emotional conditions.
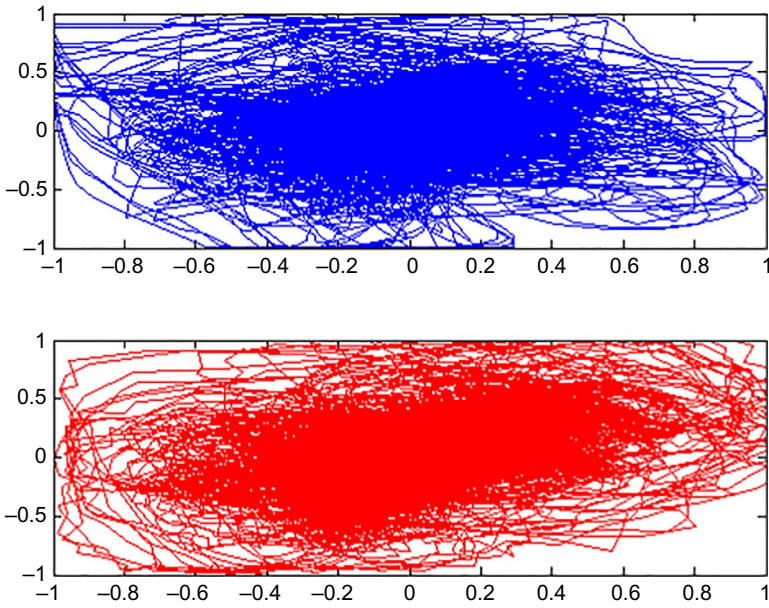
**FIG. 8.2—cont'd**

The percentage of FNN was calculated for embedding dimensions ranging from 1 to 10 with variable time delay $\tau$ varies from 1 to 60. Fig. 8.2(B) shows the fluctuation of FNN with embedding dimension for a speech signal in the angry emotion. The dimension, for which the percentage of FNN was minimum, was taken as the suitable embedding dimension.

Dimension ranging from 1 to 10 with $\tau = 43$. The phase space diagram was reconstructed for the speech signal for angry and normal emotion with proper time delay [42, 43] for each case as given in Fig. 8.2(C).

From Fig. 8.2(C), it can be seen that the number of outliers increased in the angry emotion compared to the normal emotion. It also can be seen that the normal emotion speech signal exhibited denser orbit than the angry emotion speech signal. These results suggest that different dynamics benefit different emotional states.

### 8.4.1 Noise-Free Environment

As stated above, isometry is the measure of recurrence in any phase space. It is calculated as

$$\text{Isometry} = \frac{1}{N^2} \sum_{i,j=1}^{N} R_{i,j}(\varepsilon) \tag{8.5}$$
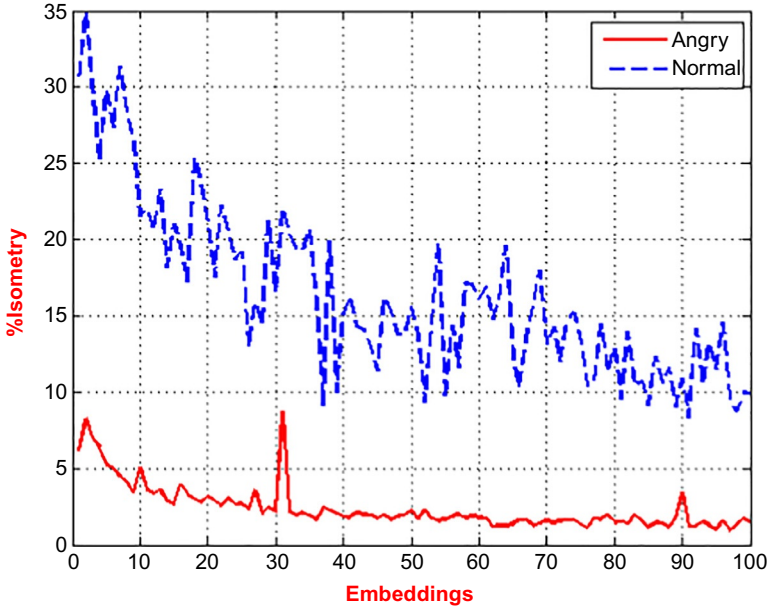
**FIG. 8.3** Embedding plot of %isometry with proper time delay for angry and normal emotional speech in a noise-free environment.

where $N$ is the number of speech samples and $R_{i,\,j}(\varepsilon)$ is the repetition matrix equivalent to a threshold of $\varepsilon$. Isometry is represented as the quantity of isometric recurrences articulated as a fraction of the entire quantity of pairs of vectors contrast in the sample (i.e., %isometry) [44]. The design of %isometry or any of its copied frameworks as a function of embedding dimension is considered as an embedding plot. The embedding plot of %isometry, represented in Fig. 8.3, shows the comparison between the %isometry of two speech signals in angry and normal emotion in a noise-free environment.

Fig. 8.3 clearly shows that the %isometry of the normal emotional speech signal was higher than that of the angry speech signal. It suggests that the autocorrelation was higher in the case of the normal speech signal. The same observation is also made for %consecutive isometry, represented in Fig. 8.4. (See Fig. 8.5.)

## 8.4.2 Noisy Environment

To investigate the effect of noise on the isometry, a Gaussian noise was added, given by $\varphi(\xi) = e^{\frac{-\xi^2}{2}}$, where $\varphi(\xi)$ is a Gaussian random variable.

When the energy of the signal is strong in the region of a restricted time interim (especially if its total energy is limited), one may calculate the energy spectral compactness. However, more commonly used is the power spectrum.
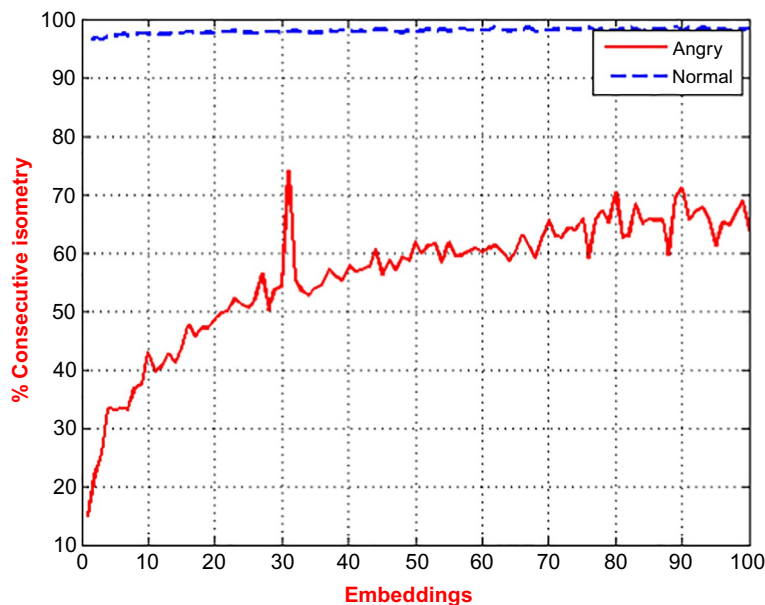
**FIG. 8.4** Embedding plot of %consecutive isometry with proper time delay for angry and normal emotional speech in a noise-free environment.
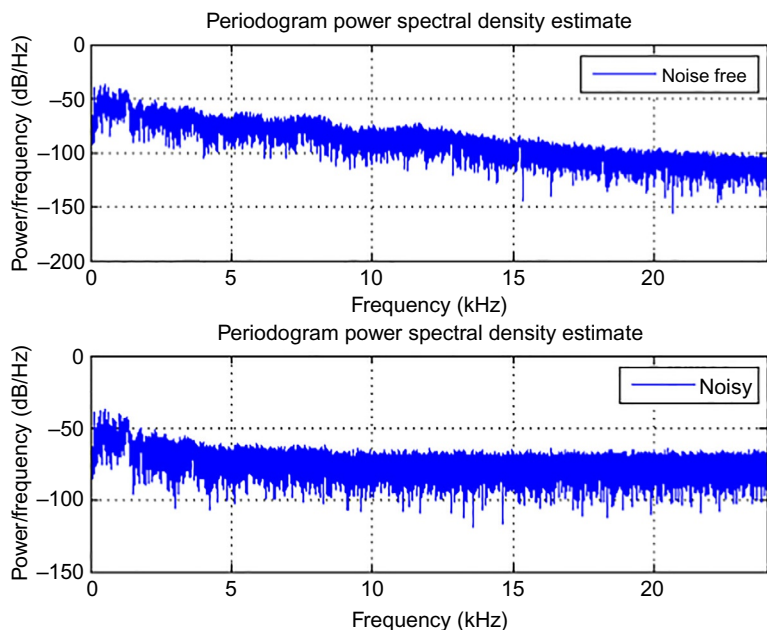


**FIG. 8.5** Power spectral density estimation (spectrum) of one speech signal in an angry emotion in a noise-free and noisy environment.
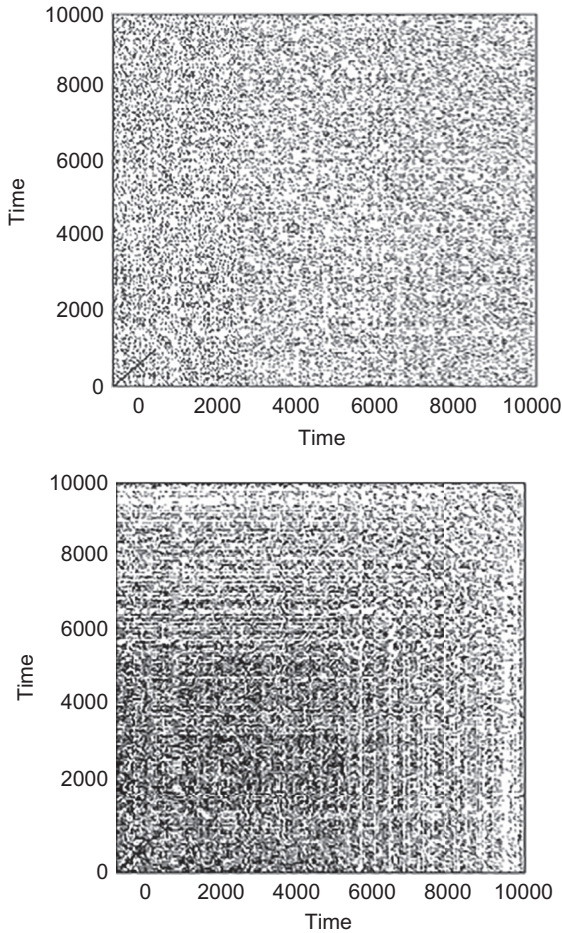
**FIG. 8.6** (A) Recurrence plot of an angry emotional speech signal with s proper time delay and embedding dimension after the addition of noise. (B) Recurrence plot of a normal emotional speech signal with a proper time delay and embedding dimension after the addition of noise.

In the experiment, it was noticed that the power frequency of the signal became higher in the power spectrum after the addition of noise.

Fig. 8.6(A) and (B) shows the recurrence plot of a speech signal in angry emotion and usual emotion, respectively. After the addition of noise.

In a noisy condition, it was observed that the recurrence plots were more complex with lots of dark dots. It suggested that the addition of noise in the signal increased the recurrence rate.

Both %isometry and %consecutive isometry were again calculated for the speech signals in angry and normal emotions after adding the noise. Figs. 8.7 and 8.8 show the embedding plots of %isometry and %consecutive isometry,
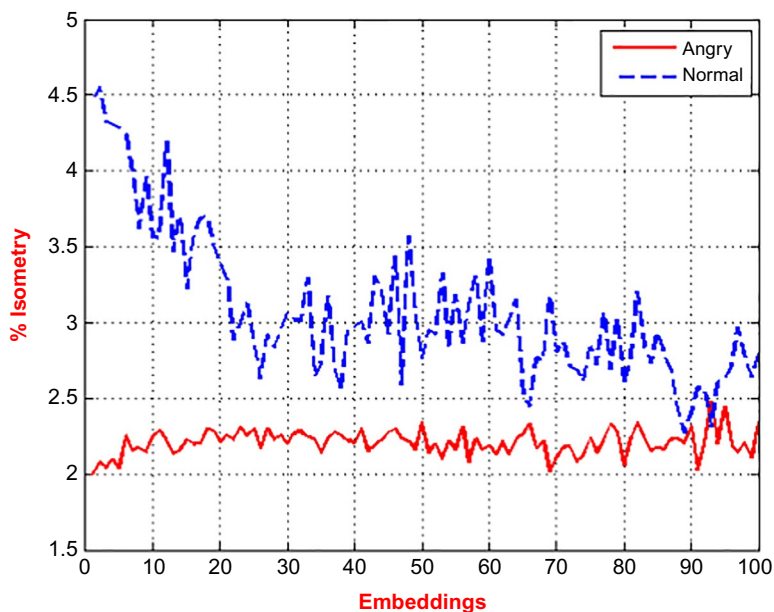
**FIG. 8.7**    Embedding plot of %isometry with a proper time delay for angry and normal emotional speech.
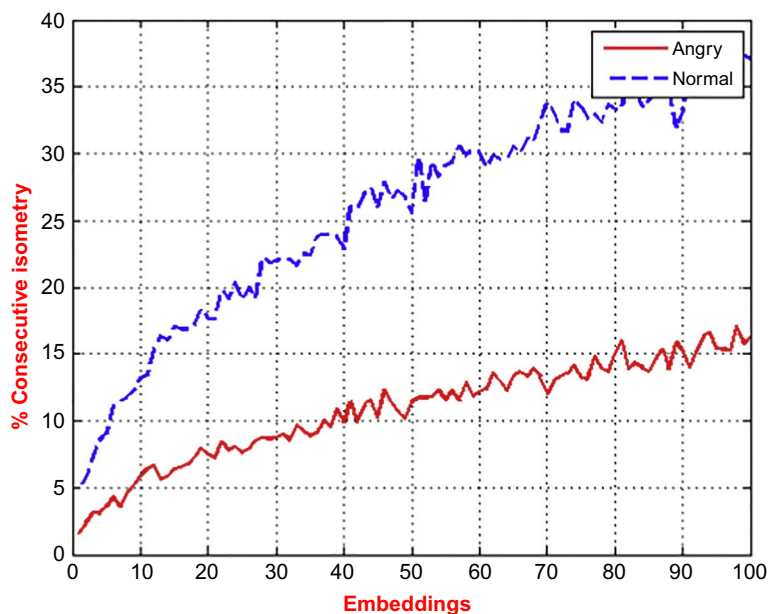


**FIG. 8.8**    Embedding plot of %consecutive isometry with a proper time delay for angry and normal emotional speech.

for speech signals in angry and normal emotion, respectively, after adding noise. It is evident from the figures that the speech signal in normal emotion exhibited more %isometry and %consecutive isometry than the speech signal in an angry emotion, similar to that of a noise-free condition. These results indicated the robustness of the two parameters even in the noisy condition.

## 8.5   Conclusion

Speech signals were acquired from a healthy male volunteer at a sampling rate of 16 KHz, when the volunteer uttered eight sentences in the Bengali language in normal and angry emotion. Phase space analysis and recurrence analysis were performed on the speech signals to classify normal and angry emotion conditions. During the editing of the phase space diagram, the proper values of time delay and the embedding dimension were approximated using auto mutual information (AMI) and false nearest neighbor (FNN) mechanism, respectively. The proper time delay was found to be more for normal emotion than angry emotion. The phase space plot for normal emotion had a denser orbit and lower number of outliers than that of the angry emotion. The recurrence plots were drawn and correlated parameters, specifically %isometry and %consecutive isometry, were calculated for the speech signals acquired during normal and angry emotion in noise-free and noisy environments. In both noise-free and noisy environments, %isometry and %consecutive isometry values were greater for normal emotion than angry emotion. This suggested that these parameters are robust to environmental noise and can be used to develop a SER system for classifying different emotional states.

## References

[1] N.J. Gogoi, J. Kalita, Emotion recognition from acted Assamese speech, Int. J. Innov. Res. Sci. Eng. Technol. 4 (6) (2015).

[2] D. Ververidis, C. Kotropoulos, Emotional speech recognition: resources, features, and methods, Speech Comm. 48 (2006) 1162–1181.

[3] M. El Ayadi, et al., Survey on speech emotion recognition: features, classification schemes, and databases, Pattern Recogn. 44 (2011) 572–587.

[4] B. Schuller, G. Rigoll, M. Lang, Hidden Markov model-based speech emotion recognition, in: 2003 International Conference on Multimedia and Expo. ICME'03. Proceedings (Cat. No.03TH8698), Baltimore, MD, USA, 2003, pp. 1–401.

[5] B. Schuller, G. Rigoll, M. Lang, Speech emotion recognition combining acoustic features and linguistic information in a hybrid support vector machine-belief network architecture, in: 2004 IEEE International Conference on Acoustics, Speech, and Signal Processing, Montreal, Que, 2004, pp. 1–577.

[6] Y.-L. Lin, G. Wei, Speech emotion recognition based on HMM and SVM, in: 2005 International Conference on Machine Learning and Cybernetics, Guangzhou, China, vol. 8, 2005, pp. 4898–4901.

[7] S. Lalitha, et al., Speech emotion recognition, in: 2014 International Conference on Advances in Electronics, Computers and Communications (ICAECC), 2014, pp. 1–4.

[8] M.S. Kamal, L. Chowdhury, M.I. Khan, A.S. Ashour, J.M.R. Tavares, N. Dey, Hidden Markov model and Chapman Kolmogrov for protein structures prediction from images, Comput. Biol. Chem. 68 (2017) 231–244.

[9] N. Dey, A.S. Ashour, W.S. Mohamed, N.G. Nguyen, Acoustic sensors in biomedical applications, in: Acoustic Sensors for Biomedical Applications, Springer, Cham, 2019, pp. 43–47.

[10] N. Dey, A.S. Ashour, W.S. Mohamed, N.G. Nguyen, Acoustic wave technology, in: Acoustic Sensors for Biomedical Applications, Springer, Cham, 2019, pp. 21–31.

[11] S.G. Koolagudi, et al., IITKGP-SESC: speech database for emotion analysis, in: International Conference on Contemporary Computing, 2009, pp. 485–492.

[12] R. Altrov, H. Pajupuu, The influence of language and culture on the understanding of vocal emotions, J. Estonian Finno-Ugric Linguistics 6 (3) (2015).

[13] R. Cowie, N. Sussman, A. Ben-Ze'ev, Emotions: concepts and definitions, in: P. Petta, C. Pelachaud, R. Cowie (Eds.), Emotionoriented Systems: The HUMAINE Handbook, Springer, Berlin, Heidelberg, 2011, pp. 9–31.

[14] E. Douglas-Cowie, N. Campbell, R. Cowie, P. Roach, Emotional speech: towards a new generation of databases, Speech Comm. 40 (1) (2003) 33–60.

[15] C.-H. Park, K.-B. Sim, Emotion recognition and acoustic analysis from speech signal, in: Proceedings of the International Joint Conference on Neural Networks, 2003, Portland, OR, vol. 4, 2003, pp. 2594–2598.

[16] S. Ntalampiras, N. Fakotakis, Modeling the temporal evolution of acoustic parameters for speech emotion recognition, IEEE Trans. Affect. Comput. 3 (1) (2012) 116–125.

[17] L. Zão, D. Cavalcante, R. Coelho, Time-frequency feature and AMS-GMM mask for acoustic emotion classification, IEEE Signal Process. Lett. 21 (5) (2014) 620–624.

[18] J. Chang, X. Zhang, Q. Zhang, Y. Sun, Investigating duration effects of emotional speech stimuli in a tonal language by using event-related potentials, IEEE Access 6 (2018).

[19] Y. Kang, A. Li, Prosody conversion from neutral speech to emotional speech, IEEE Trans Audio Speech Lang. Process. 14 (4) (2006) 1145–1154.

[20] C. Wu, C. Hsia, C. Lee, M. Lin, Hierarchical prosody conversion using regression-based clustering for emotional speech synthesis, IEEE Trans. Audio Speech Lang. Process. 18 (6) (2010) 1394–1405.

[21] J. Jia, S. Zhang, F. Meng, Y. Wang, L. Cai, Emotional audio-visual speech synthesis based on PAD, IEEE Trans. Audio Speech Lang. Process. 19 (3) (2011) 570–582.

[22] N. Dey, A.S. Ashour, Direction of Arrival Estimation and Localization of Multi-Speech Sources, Springer International Publishing, 2018.

[23] N. Dey, A.S. Ashour, Applied examples and applications of localization and tracking problem of multiple speech sources, in: Direction of Arrival Estimation and Localization of Multi-Speech Sources, Springer, Cham, 2018, pp. 35–48.

[24] E. Navas, I. Hernaez, I. Luengo, An objective and subjective study of the role of semantics and prosodic features in building corpora for emotional TTS, IEEE Trans. Audio Speech Lang. Process. 14 (4) (2006) 1117–1127.

[25] S. Zhang, S. Zhang, T. Huang, W. Gao, Speech emotion recognition using deep convolutional neural network and discriminant temporal pyramid matching, IEEE Trans. Multimedia 20 (6) (2018) 1576–1590.

[26] F. Weninger, F. Eyben, B. Schuller, On-line continuous-time music mood regression with deep recurrent neural networks, in: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Florence, 2014, 2014, pp. 5412–5416.

[27] A.R. Avila, J. Monteiro, D. O'Shaughneussy, T.H. Falk, Speech emotion recognition on mobile devices based on modulation spectral feature pooling and deep neural networks, in: IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), Bilbao, 2017, 2017, pp. 360–365.

[28] N. Dey, A.S. Ashour, Challenges and future perspectives in speech-sources direction of arrival estimation and localization, in: Direction of Arrival Estimation and Localization of Multi-Speech Sources, Springer, Cham, 2018, pp. 49–52.

[29] N. Dey, A.S. Ashour, Sources localization and DOAE techniques of moving multiple sources, in: Direction of Arrival Estimation and Localization of Multi-Speech Sources, Springer, Cham, 2018, pp. 23–34.

[30] N. Dey, A.S. Ashour, Microphone array principles, in: Direction of Arrival Estimation and Localization of Multi-Speech Sources, Springer, Cham, 2018, pp. 5–22.

[31] Y. Zong, W. Zheng, T. Zhang, X. Huang, Cross-corpus speech emotion recognition based on domain-adaptive least-squares regression, IEEE Signal Process. Lett. 23 (5) (2016) 585–589.

[32] B. Sivakumar, et al., River flow forecasting: use of phase-space reconstruction and artificial neural networks approaches, J. Hydrol. 265 (2002) 225–245.

[33] S.P. Chandrasekaran, A nonlinear dynamic modelling for speech recognition using recurrence plot—a dynamic Bayesian approach, in: IEEE International Conference on Signal Processing and Communications, ICSPC 2007, 2007, pp. 516–519.

[34] A.M. Fraser, H.L. Swinney, Independent coordinates for strange attractors from mutual information, Phys. Rev. A 33 (1986) 1134.

[35] A. Dey, et al., A new kind of dynamical pattern towards distinction of pre-meditative and meditative states through HRV, Science 3 (2012).

[36] M.B. Kennel, H.D. Abarbanel, False neighbors and false strands: a reliable minimum embedding dimension algorithm, Phys. Rev. E 66 (2002) 026209.

[37] A. Dey, D.K. Bhattacha, D.N. Tibarewala, N. Dey, A.S. Ashour, D.-N. Le, E. Gospodinova, M. Gospodinov, Chinese-chi and kundalini yoga meditations effects on the autonomic nervous system: comparative study, Int. J. Interact. Multimedia Artif. Intell. 3 (7) (2016) 87–95. 9 p.

[38] M.B. Kennel, et al., Determining embedding dimension for phase-space reconstruction using a geometrical construction, Phys. Rev. A 45 (1992) 3403.

[39] J.P. Eckmann, et al., Recurrence plots of dynamical systems, EPL (Europhys. Lett.) 4 (1987) 973.

[40] N. Marwan, et al., Recurrence plots for the analysis of complex systems, Phys. Rep. 438 (2007) 237–329.

[41] S.K. Nayak, K. Pande, P.K. Patnaik, S. Nayak, S.J. Patel, A. Anis, A. Dey, K. Pal, Understanding the effect of cannabis abuse on the ANS and cardiac physiology of the Indian women paddy-field workers using RR interval and ECG signal analyses. in: IEEE International Conference APSIPA ASC 2017, Aloft Kuala Lumpur Sentral Sentral, Kuala Lumpur, 2017. https://doi.org/10.1109/APSIPA.2017.8282047.

[42] A. Dey, S.K. Palit, S. Mukherjee, D.K. Bhattacharya, D.N. Tibarewala, A new technique for the classification of pre-meditative and meditative states, in: IEEE International Conference, "ICCIA-2011", 2011. Print ISBN 978-1-4577-1915-8.

[43] M. Das, T. Jana, P. Dutta, R. Banerjee, A. Dey, D.K. Bhattacharya, M.R. Kanjilal, Study the effect of music on HRV signal using 3D poincare plot in spherical co-ordinates—a signal processing approach, in: IEEE International Conference on Communication and Signal Processing, April 2–4, 2015, India, 2015. ISBN 978-1-4799-8080-2.

[44] H. Sabelli, A. Lawandow, Homeobios: the pattern of heartbeats in newborns, adults, and elderly patients, in: Nonlinear Dynamics, Psychology, and Life Sciences, vol. 14, 2010, p. 381.