

Dialect Identification of Assamese Language using Spectral Features

Tanvira Ismail* and L. Joyprakash Singh

Department of Electronics and Communication, School of Technology, North-Eastern Hill University, NEHU Campus, Shillong – 793022, Meghalaya, India; tanvira.ismail@gmail.com

Abstract

Objective: Accurate dialect identification technique helps in improving the speech recognition systems that exist in most of the present day electronic devices and is also expected to help in providing new services in the field of e-health and telemedicine which is especially important for older and homebound people. **Methods:** In this paper we have developed the speech corpora needed for the dialect identification purpose and described a method to identify Assamese language (an Indian language) and two of its dialects, viz., Kamrupi and Goalparia. **Findings:** Research work done on dialect identification is relatively much less than that on language identification for which one of the reasons being dearth of sufficient database on dialects. As mentioned, we have developed the database and then Mel-Frequency Cepstral Coefficient has been used to extract the spectral features of the collected speech data. Two modeling techniques, namely, Gaussian Mixture Model and Gaussian Mixture Model with Universal Background Model have been used as the modeling techniques to identify the dialects and language. **Novelty:** So far, standard speech database for Assamese dialects that can be used for speech processing research has not been made. In this paper, we not only describe a method to identify dialects of the Assamese language, but have also developed the speech corpora needed for the dialect identification purpose.

Keywords: Assamese, Gaussian Mixture Model, Gaussian Mixture Model with Universal Background Model, Goalparia, Kamrupi, Mel-Frequency Cepstral Coefficient

1. Introduction

Dialect is defined as a variety of a language spoken in a particular geographical area that is distinguished from other varieties of the same language by features of pronunciation, grammar and vocabulary. It is also defined as a variety of speech that differs from the standard language normally designated as the official language.

Humans are born with the ability to discriminate between spoken languages as part of human intelligence and the quest to automate such ability has never stopped.¹ Just like any other artificial intelligence technologies, automatic dialect identification aims to replicate such human ability through computational means.¹ Dialect identification is the task of recognizing a speaker's regional dialect within a predetermined language.

Developing a good method to detect dialect accurately helps in

- Improving certain applications and services such as Speech Recognition Systems which exist in most of the today's electronic devices,
- Enhancing the human - computer interaction applications and
- Securing the remote access communication.²

Moreover, accurate dialect detection technique is expected to help in providing new services in the field of e-health and telemedicine which is especially important for older and homebound people.²

Although much research work has been done in language identification, the problem of dialect identification, which is very similar to the problem of language identification, has not received the same level of research interest.³ Research in the field of dialect is still limited due to the dearth of databases and the time consuming analysis process.⁴ Initial work using Gaussian Mixture

*Author for correspondence

Models (GMM) was performed for the identification of twelve languages and dialect identification was also performed for three out of the twelve languages, viz., English, Mandarin and Spanish.³ Identification capability was improved by using Gaussian Mixture Model with Universal Background Model (GMM-UBM) and it was showed that a GMM-UBM based model provides good results for recognition of American vs. Indian English, four Chinese dialects and three Arabic dialects.⁵

In the field of Indian languages, Mel-Frequency Cepstral Coefficient (MFCC) and Speech Signal based Frequency Cepstral Coefficient (SFCC) feature extraction techniques along with GMM and Support Vector Machine (SVM) modeling techniques was used to identify the two main language families of India, viz., Indo-Aryan and Dravidian, which in total consisted of 22 languages.⁶ The Indo-Aryan family consisted of 18 languages, viz., Assamese, Bengali, Bhojpuri, Chhattisgarhi, Dogri, English, Gujarati, Hindi, Kashmiri, Konkani, Manipuri, Marathi, Nagamese, Odia, Punjabi, Sanskrit, Sindhi and Urdu, while Dravidian family consisted of 4 languages, viz., Kannada, Malayalam, Tamil and Telugu.⁶ A language identification system with two levels that used acoustic features, was modeled using Hidden Markov Model (HMM), GMM and Artificial Neural Network (ANN) and was tested on nine Indian languages, viz., Tamil, Telugu, Kannada, Malayalam, Hindi, Bengali, Marathi, Gujarati and Punjabi.⁷ In this two level language identification system, the family of the spoken language is identified in the first level and after feeding this input to the second level the identification of a particular language is made in the corresponding family.⁷ On the other hand, both spectral features and prosodic features were used for analyzing the specific information with regards to each language present in speech and then GMM was applied for identification purposes on the Indian language speech database (IITKGP-MLILSC) which consists of as many as 27 Indian languages.⁸ Concentrating more on South Indian languages (languages of Dravidian family), a Language Identification (LID) system was presented that worked for the four South Indian languages, viz., Kannada, Malayalam, Tamil and Telugu and one North Indian language, viz., Hindi in which each language was modeled using an approach based on Vector Quantization, whereas the speech was segmented into dierrant sounds and the performance of the system on each of the segments was studied.⁹ An LID system using GMM for the features that were extracted was further modeled using

Split and Merge Expectation Maximization Algorithm was tested on four Indian languages, viz, Hindi, Telugu, Gujarati and English.¹⁰ Four Indian languages, viz, Hindi, Bengali, Oriya and Telugu were identified by considering the special CV (Consonant-Vowel) behavior of the language in their syllables and were also analyzed using the SVM classifier.¹¹ Work was also done on language identification that was speaker dependent and was tested on three Indian languages, viz., Assamese, Hindi and Indian English, based on clustering and supervised learning.¹² First the feature vectors using LP coefficients were obtained and clusters of vectors using the K-means algorithm were formed. Supervised learning was then used for recognizing the probable cluster to which the test speech sample belongs.¹²

As far as Indian dialects are concerned, both spectral features and prosodic features were used to identify five Hindi dialects, viz., Chattisgarhi, Bengali, Marathi, General and Telugu using Autoassociative Neural Network (AANN) models.¹³ For each dialect, their database consisted of data from 10 speakers speaking in spontaneous speech for about 5-10 minutes resulting in a total of 1-1.5 hours.¹³ On the other hand, speaker identification of specific dialect, viz., Assamese was done using features obtained from various speaker dependent parameters of voiced speech and then ANN based classifiers were used for identification purpose.¹⁴

Thus, with regard to Indian languages vis-à-vis Indian dialects also, the dialects have not been explored as much as Indian languages. Therefore, the present focus of our research is on Indian dialects identification and in this paper we have chosen for our study the dialects of Assamese language. Assamese has mainly three dialects, viz., the standard dialect, Kamrupi and Goalparia dialects.¹⁵ The standard dialect is the standard Assamese language which is an Indo-Aryan language designated as the official language of Assam state located in the north-eastern part of India. In Assam, Assamese is the formal written language used for education, media and other official purposes, while Kamrupi and Goalparia are the informal spoken dialects spoken in the Kamrup and Goalparia districts of Assam, respectively. The importance of Kamrupi dialect can be understood from the early literature, wherein Kamrupi was documented as the first ancient Aryan literary language spoken both in the Brahmaputra valley of Assam and North Bengal.¹⁶⁻¹⁸ On the other hand, the importance of Goalparia, which is also an Indo-Aryan dialect, is evident from the fact that it has

a rich culture of folk music known as Goalparia Lokgeet. According to the 2011 Census of India, there are about 20, 6 and 1 million speakers of Assamese language, Kamrupi dialect and Goalparia dialect, respectively.

So far, standard speech database for Assamese dialects that can be used for speech processing research has not been made. In fact, as mentioned above, dearth of databases is one of the main reasons for limited research work on dialect identification. Therefore, in this paper, we not only describe a method to identify dialects of the Assamese language, but have also developed the speech corpora needed for the dialect identification purpose. In Section 2 of the paper various aspects of the experimental setup and the results obtained are described. The accuracy of the system is discussed in section 3 and conclusion is given in Section 5.

2. Experimental Setup and Results Obtained

Any identification system must consist of two phases i.e. training and testing. The training phase is used to build the system while the testing phase is used to check whether the system built during the training phase is working properly. Figure 1 gives an idea about the training and testing phases of the present identification system.

2.1 Development of Speech Corpora

While building a speech corpus, the important criteria to be kept in mind are:

- Enough speech must be recorded from enough speakers in order to validate an experiment under study,¹⁹
- Speakers of wide ranging age must be considered to study the effects of age on pronunciation,²⁰
- Speakers of varying educational backgrounds must be considered in order to track the effects of speech effort level and speaking rate, intelligence level and speaker's experience,^{21,22}
- Recording of data must be done in different ranges of environments such as home, office, roads, cars, noisy conditions or villages, and²⁰
- Data must also be collected with multiple training and testing sessions in order to track the effects of intersession variability on the identification system.^{19,23}

A speech corpus can be of two types: read speech or spontaneous speech.²⁴ In read speech, people are asked to read a written text while in spontaneous speech; there can be narratives where a single person is asked to speak on a topic of his choice.²⁴

For the present work, speech data has been recorded using the Zoom H4N Handy Portable Digital Recorder

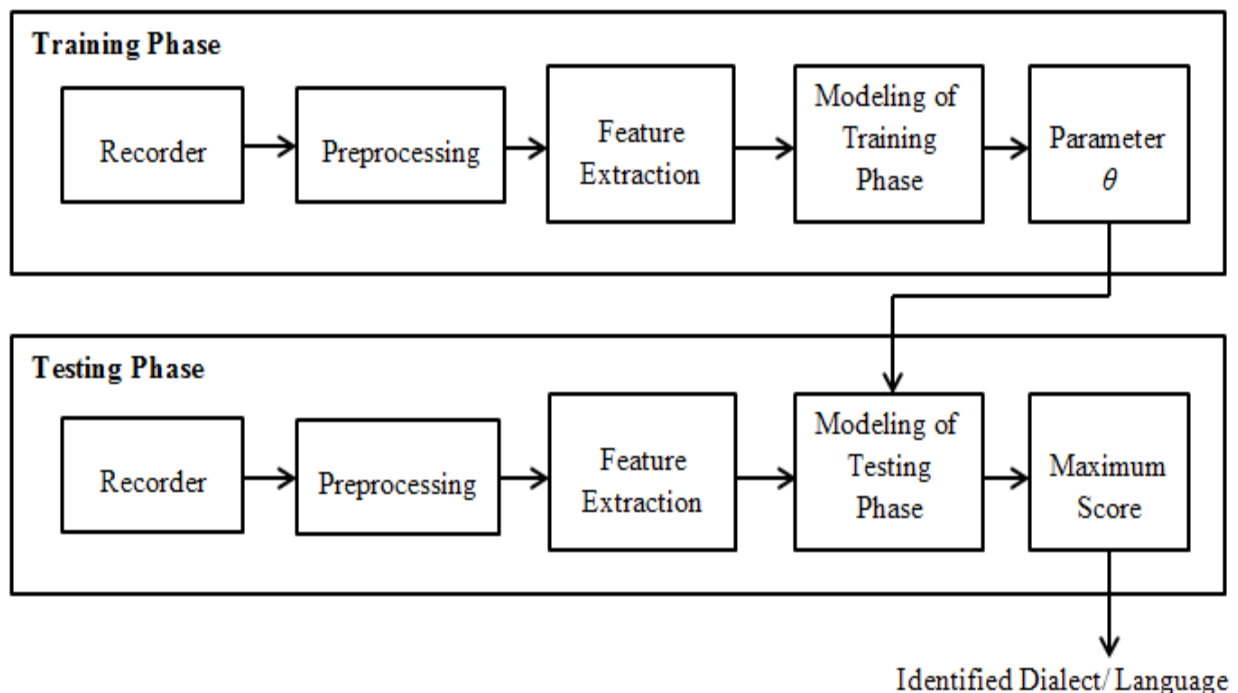


Figure 1. Block diagram of the present identification system.

and 8 kHz sampling frequency. Recording has been done as spontaneous speech. Speakers mostly spoke about their childhood, home town, career, personal habits and so on. Rather than reading a study material or speaking fixed text sentences, allowing the speakers to speak continuously on topics of their choice, maintains the speaker’s natural manner of speaking. Data has been collected from people residing in both cities and villages coming from different walks of life such as college students, teachers, office goers and farmers. Data have been recorded at homes or offices or classrooms situated both in cities and villages by taking care to avoid noisy conditions during recording.

Speech data of 6 hours and 8 minutes from 38 speakers for Assamese language, of 4 hours and 22 minutes from 30 speakers for Kamrupi dialect and of 3 hours from 27 speakers for Goalparia dialect have been collected. The speakers chosen were in the age group of 21-60 years.

During preprocessing, the recorded speech data have been listened to and analyzed carefully. It has been observed that people from the city very often make use of English words even when conversing in their native dialect or language. Care has been taken to remove such words and any other portions that are not part of the present dialects or language. The wave files have also been cut into smaller portions of 3 to 7 seconds each. The description of the corpora made for the present identification system is also summarized in Table 1.

2.2 Feature Extraction

The time domain waveform of a speech signal carries all possible acoustic information and from the point of view of phonology, very little can be said on the basis of the waveform itself.²⁵ To be able to find some relevant information from incoming data, it is important to have some methods to reduce the information present in each segment of the audio signal into a comparatively small number of parameters or features.²⁵ Feature extraction is defined as the process of conserving the important information present in the speech signal while removing the unwanted portions. The spectral features of a speech

signal are obtained by converting the time domain signal into the frequency domain. To obtain the spectral features, MFCC has been used as the feature extraction technique for both the training phase and testing phase of the present identification system. In the present system, the speech signal is segmented into frames of 20 milliseconds with an overlap of 10 milliseconds and each frame is then multiplied by a hamming window. After windowing, Fast Fourier Transform (FFT) is performed in order to acquire the magnitude frequency response of each frame. The magnitude frequency response is then multiplied by triangular band pass filters to get the log energy of each triangular band pass filter. In the present system, 22 triangular band pass filters have been used. The relationship between the common linear frequency of speech (f) and the Mel-frequency is given in Equation 1.

$$\text{Mel}(f) = 1125 \ln \left(1 + \frac{f}{700} \right) \tag{1}$$

Discrete Cosine Transform (DCT) is then applied to the 22 log energies acquired from the triangular band pass filters in order to get the MFCC. The expression for DCT is given in Equation 2.

$$C_m = \sum_{k=1}^N \cos[m(k-0.5)\pi / N] E_k ; m = 1,2,\dots,L \tag{2}$$

Where N is the number of triangular band pass filters, L is the number of MFCCs and E_k is the log energies acquired from the triangular band pass filters. In the present system, N = 22 and L = 14. The MFCCs obtained are then used as an input to the modeling techniques.

2.3 Modeling Technique

Two modeling techniques have been used, viz. GMM and GMM-UBM. Mixture models are a type of density model which comprise of a number of density functions, usually Gaussian and the component functions are combined to provide a multimodal density.²⁶ Gaussian belongs to the exponential family. This family has same spectrum in both frequency and time domain. GMM is defined as

Table 1. Statistical description of the speech corpora (Recording environment: Homes / offices / classrooms in cities and villages)

Language / Dialect	Total duration of speech data	Number of Speakers	Age group of Speakers	Sampling frequency	Duration of wave file
Assamese	6 h 8 min	38	21 – 60 yr	8 kHz	3 – 5 s
Kamrupi	4 h 22 min	30	21 – 60 yr	8 kHz	4 – 7 s
Goalparia	3 h	27	21 – 60 yr	8 kHz	4 – 7 s

a parametric Probability Density Function (PDF) that is represented as a weighted sum of Gaussian component density. Parameters of the Gaussian Mixture Model consist of the mean and variance matrices of the Gaussian components and the weights indicating the contribution of each Gaussian to the approximation of PDF.

The MFCCs of training data obtained after feature extraction are fed into the training phases of GMM and GMM-UBM. Once the training phase is over, the mean, variance and weight, which are collectively known as parameter θ , are obtained for GMM. A separate GMM has to be built for each of the constituent dialects and languages. For the present system, three GMMs are built one each for Assamese language, Kamrupi dialect and Goalparia dialect.

Unlike GMM only one UBM is built to obtain θ for all constituent dialects and languages. Hence, the name Universal is appropriately given since it is trained with data from different dialects and languages. In the present identification system, initially, the MFCCs of Assamese language, Kamrupi and Goalparia dialects are used to train just one GMM which is called the Universal Background Model (UBM). This initial θ contains information of all the dialects and languages being used in the system. For the present system, this initial θ contains information of Assamese language, Kamrupi and Goalparia dialects. From the initial UBM, a distinct model for every constituent dialects and languages is derived by MAP (Maximum A-Posteriori) adapting the trained GMM-UBM to the training data of that dialect or language. This gives the adapted parameters of the constituent dialects and languages. For the present identification system, after MAP adaptation, the adapted parameter θ is obtained for Assamese language, Kamrupi and Goalparia dialects.

2.4 Results Obtained

The speech signals used in the testing phase have to be different from those in the training phase. Also, the speech signals used should be a combination of all the languages and dialects for which the GMMs and GMM-UBM have been built during the training phase. In the present system, the speech signals used in the testing phase consist of Assamese language, Kamrupi and Goalparia dialects. The testing phase of GMM is loaded with the parameter θ while the testing phase of GMM-UBM is loaded with the adapted parameter θ obtained from training phases of GMM and GMM-UBM, respectively. In the present system, three θ parameters and three adapted θ parameters are loaded, one each for Assamese language, Kamrupi and Goalparia dialects.

A language or dialect is detected when a test utterance gets maximum likelihood score for that particular language or dialect. For the present system, when a test utterance gets maximum likelihood score for Kamrupi dialect, it implies that, the test utterance is spoken in Kamrupi. The same holds true for Assamese language and Goalparia dialect. In Figure 2, snapshots of some of the likelihood scores obtained in the present identification system have been shown. Column 1 stands for Assamese language, 2 for Goalparia dialect and 3 for Kamrupi dialect. The rows are for test utterances. In Figure 2(a), all test utterances except the first test utterance, are getting maximum likelihood score in Column 1 which means these test utterances are spoken in Assamese language while the first test utterance has been wrongly detected as Kamrupi dialect. In Figure 2(b), all test utterances are getting maximum likelihood score in Column 2 and hence are spoken in Goalparia dialect. In Figure 2(c), all test utterances are

allLiksum <1149x3 double>			
	1	2	3
1	-16.4513	-15.8105	-15.0702
2	-31.4544	-34.5559	-34.6016
3	-26.8436	-29.6805	-29.5825
4	-33.7616	-37.3927	-37.2158
5	-30.2046	-33.3158	-33.5200
6	-32.7819	-36.0649	-35.6579
7	-32.7895	-36.1852	-36.1196
8	-32.6199	-36.3940	-36.3862
9	-29.7023	-32.8582	-32.8930
10	-30.6554	-33.8248	-33.5820

(a)

allLiksum <1149x3 double>			
	1	2	3
734	-17.4550	-14.0644	-14.4079
735	-20.8434	-16.9991	-17.4952
736	-19.2586	-15.1223	-15.7088
737	-19.2593	-15.6995	-16.2088
738	-21.1341	-17.6936	-18.1429
739	-18.4926	-14.2206	-15.0158
740	-19.2961	-14.6964	-15.5431
741	-17.3450	-12.1212	-13.7539
742	-19.4247	-14.5928	-15.9152
743	0.7358	5.3451	2.1826

(b)

allLiksum <1149x3 double>			
	1	2	3
1126	-23.1415	-22.0786	-20.4122
1127	-22.5678	-21.1308	-19.6947
1128	-20.6087	-19.3397	-17.9204
1129	-21.0884	-20.2123	-18.4608
1130	-21.8622	-20.3770	-18.7929
1131	-21.4275	-19.9659	-18.5484
1132	-23.1707	-21.5475	-19.8894
1133	-24.1147	-23.1439	-21.3544
1134	-22.1137	-20.8754	-19.1270
1135	-21.9410	-20.9085	-19.1955

(c)

Figure 2. Some of the resulting likelihood scores of the present system. (a) Likelihood scores where the test utterances have been detected as Assamese. (b) Likelihood scores where the test utterances have been detected as Goalparia. (c) Likelihood scores where the test utterances have been detected as Kamrupi.

Table 2. Statistical description of training and testing data used

Language / Dialect	Training data used	Testing data used	Length of wave files	No. of testing files
Assamese	4.16 h	33 mins	3-5 s	395
Kamrupi	3.73 h	38 mins	4-7 s	389
Goalparia	2.3 h	31 mins	4-7 s	365

Table 3. Calculation of accuracy of the identification system

Language / Dialect	No. of Testing Files used	Total No. of Testing Files	Total No. of False Identification		Accuracy	
			GMM	GMM-UBM	GMM	GMM-UBM
Assamese	395	1149	164	19	85%	98.3%
Kamrupi	389					
Goalparia	365					

getting maximum likelihood score in column 3 and hence are spoken in Kamrupi dialect.

The problem of dialect identification is viewed more challenging than that of language recognition due to the greater similarity between dialects of the same language.²⁷ In the present system too, the wrong detections are due to the excessive similarities between Assamese language, Kamrupi and Goalparia dialect.

3. Calculation of Accuracy

For Assamese language, the training data set used is of 4.16 hours and testing data set used is of 33 minutes. A total of 395 testing files have been used. For Kamrupi dialect, the training data set used is of 3.73 hours and testing data set used is of 38 minutes. A total of 389 testing files have been used. For Goalparia dialect, the training data set used is of 2.3 hours and testing data set used is of 31 minutes. A total of 365 testing files have been used. As mentioned above, each wave file is of 3 to 7 seconds long. All these aspects of the training and testing data used for the present identification system are also summarized in Table 2.

The performance of a dialect identification system is determined by the Identification Rate (IDR). The IDR is calculated by using Equation 3.

$$IDR = \frac{n}{N} \quad (3)$$

Where n is the number of correctly identified utterances and N is the total number of utterances. For the present identification system, a total of 1149 testing files

(395 Assamese, 389 Kamrupi and 365 Goalparia) have been used and out of these false detections for GMM and GMM-UBM were 164 and 19, respectively. The present system is thus giving an accuracy of 85.7% for GMM and 98.3% for GMM-UBM. These results are summarized in Table 3.

4. Conclusion

In the present work, a speech corpus for two dialects and one language has been built with a total duration of 13 hours and 30 minutes using spontaneous speech. Although Assamese language, Kamrupi and Goalparia dialects are very similar to each other, the present work shows that both GMM and GMM-UBM can be successfully used to identify them. When GMM is used, the system identifies the constituent dialects and language with an accuracy of 85.7%, whereas the system gives an accuracy of 98.3% when GMM-UBM is used. Thus, the present identification system gives a better accuracy with GMM-UBM as the modeling technique.

For future work, the database can be increased to see if it changes the accuracy of the system. Moreover, for future work prosodic features can be used to see if it improves the identification capability of the system, since in the present system spectral features have been used.

5. References

1. Li H, Ma B, Lee AK. Spoken language recognition: From fundamentals to practice. Proceedings of the IEEE. 2013; 101(5):1136–59. Crossref

2. Etman A, Louis BAA. Language and dialect identification: A survey. Proceedings of SAI Intelligent Systems Conference; 2015. p. 220–31. Crossref
3. Torres-Carrasquillo PA, Gleason TP, Reynolds DA. Dialect identification using Gaussian Mixture Models Proceedings of the Speaker and Language Recognition Workshop; 2004. p. 41–4.
4. Chan FN, Shan W, Campbell JP. A linguistically-informative approach to dialect recognition using Dialect-Discriminating Context-Dependent Phonetic Models. Proceedings of the IEEE International Conference on Acoustic Speech Signal Processing; 2010. p. 5014–7.
5. Torres-Carrasquillo PA, Sturim D, Reynolds D, McCree A. Eigen-Channel Compensation and Discriminatively Trained Gaussian Mixture Models for dialect and accent recognition. Proceedings of the INTERSPEECH; 2008. p. 723–6.
6. Sengupta D, Saha G. Identification of the major language families of India and evaluation of their mutual influence. Current Science. 2016; 110(4):667–81. Crossref
7. Jothilakshmi S, Ramalingam V, Palanivil S. A hierarchical language identification system for Indian languages. Digital Signal Processing. 2012; 22(3):544–53. Crossref
8. Reddy VR, Maity S, Rao KS. Identification of Indian languages using Multi-Level Spectral and Prosodic Features. International Journal of Speech Technology. 2013; 16(4):489–511. Crossref
9. Ballela J, Murthy HA, Nagarajan T. Language identification from short segments of speech Proceedings of the INTERSPEECH; 2000. p. 1033–6.
10. Manwani N, Mitra SK, Joshi MV. Spoken language identification for Indian languages using split and merge EM Algorithm. Ghosh A, De RK, Pal SK. Pattern Recognition and Machine Intelligence. Editions. PReMI 2007. Lecture Notes in Computer Science. 2007; 4815:463–8. Crossref
11. Mohanty S. Phonotactic model for spoken language identification in Indian language perspective. International Journal of Computer Applications. 2011; 19(9):18–24. Crossref
12. Roy P. Language recognition of three Indian languages based on clustering and supervised learning. Proceedings of the International Conference on Computer Applications-Telecommunications; 2010. p. 77–82. Crossref
13. Rao KS, Koolagudi SG. Identification of Hindi dialects and emotions using spectral and prosodic features of speech. Journal of Systemics Cybernetics and Informatics. 2011; 9:24–33.
14. Dutta M, Patgiri C, Sarma M, Sarma K. Closed-set text-independent speaker identification system using multiple ANN classifiers. Proceedings of the 3rd International Conference on Frontiers of Intelligent Computing: Theory and Applications; 2014. p. 377–85. PMID:25471364
15. Language Information Service (LIS) – India.
16. Goswami U. A study on Kamrupi: A dialect of Assamese. 1st edition. Department of Historical Antiquarian Studies. Government of Assam: Gauhati; 1970.
17. Goswami U. An introduction to Assamese. Mani-Manik Prakash: Gauhati; 1978.
18. Medhi K. Assamese grammar and origin of the Assamese language. 1st edition. Publication Board. Government of Assam: Guwahati; 1988.
19. Reynolds DA. An overview of automatic speaker recognition technology. Proceedings of International Conference on Acoustics Speech and Signal Processing. 2002; 4:4072–5. Crossref, Crossref
20. Patil HA, Basu TK. Development of speech corpora for speaker recognition research and evaluation in Indian languages. International Journal of Speech Technology. 2008; 11:17–32. Crossref
21. Kersta LG. Voiceprint classification. Journal of the Acoustical Society of America. 1965; 37:1217. Crossref
22. Junqua JC, Fincke SC, Field K. The Lombard effect: A reflex to better communicate with other in noise. Proceedings of the IEEE International Conference on Acoustic Speech Signal Processing; 1999. p. 2083–6. Crossref
23. Gish H, Schmidt M. Text-independent speaker identification. IEEE Signal Processing Magazine. 1994; 11(4):18–32. Crossref
24. Colleen Richey. Speech Corpora; 2000.
25. Shrawanka U, Thakare V. Techniques of feature extraction in speech recognition system: A comparative study. International Journal of Computer Applications in Engineering Technology and Sciences. 2010; 2:412–8.
26. Nour-Eddine L, Abdelkader A. GMM-based Maghreb dialect identification system. Journal of Information Processing Systems. 2015; 11(1):22–38.
27. Muthusamy YK, Jain N, Cole RA. Perceptual benchmarks for automatic language identification. Proceedings of the IEEE International Conference on Acoustic Speech Signal Processing; 1994. p. 333–6. Crossref