

Acoustic Features based Accent Classification of Kashmiri Language using Deep Learning

Shehzen Sidiq Malla

Received: 5 December 2021 Accepted: 29 December 2021 Published: 9 January 2022

Abstract

Automatic identification of accents is important in today's world, where we are surrounded by ASR systems. Accent classification is the problem of knowing the native place of a person from the way He/She speaks the language into consideration. Accents are present in almost all the languages and it forms an important part of the language. Accents are produced from prosodic and articulation characteristics; in this research the aim is to classify accents of Kashmir Language. We have considered using the MFCC and Mel spectrograms for our research. A lot of research has been done for languages like English and is being done in this field and many models of machine learning and deep learning have shown state of the art results, but this problem is new for Kashmiri Language. The accents in Kashmir, vary from area to area and we have chosen 6 areas as our classes. We extracted the features from the audio data, converted those features into Images and then used the CNN architectures as our model. This research can be taken as base research for further researches in this language. Our custom models achieved the loss of 0.12 and accuracy of 98.66

Index terms— accent classification, CNN, RELU, mel-spectrograms, MFCC.

1 Introduction

ashmiri or Koshur is a Dardic language subgroup from Indo-Aryan, spoken by over seven million Kashmiris [1]. There are many accents Spoken in Kashmir. There are some major accents and some minor accents in this language. This leads to diversity in the language and adds to its beautiful sounds and variations. The aim of this research is to classify these different accents. Although many accents are being spoken in this language, for this research, we have classified the prominent accents belonging to Kupwara, Srinagar, Islamabad, Shopian, and Bandipora. The proposed approach is on the basis of using Convolution Neural Networks (CNN) and training Neural networks on the images of features extracted from the audio files. The features are Mel-spectrogram and MFCCs. Our approach uses CNN as the classifier and MFCCs, Mel spectrograms as Features. Three types of MFCCs are extracted, 13, 24 and 36. We got excellent results on our dataset.

Accent classification refers to the problem of inferring the native language of a speaker from his or her foreign accented speech. Identifying idiosyncratic differences in speech production is important for

Author: e-mail: mallashehzen786@gmail.com improving the robustness of existing speech analysis systems. For example, automatic speech recognition (ASR) systems exhibit lower performance when evaluated on foreign accented speech. By developing pre-processing algorithms that identify the accent, these systems can be modified to customize the recognition algorithm to the particular accent [1] [2]. In addition to ASR applications, accent identification is also useful for forensic speaker profiling by identifying the speaker's regional origin and ethnicity in applications involving targeted marketing [3] [4]. In this paper we propose a method for classification of 11 accents directly from the speech acoustics.

For example, Deshpande et al. used GMMs based on formant frequency features to discriminate between standard American English and Indian accented English [6]. Chen et al. explored the effect of the number of components in GMMs on classification performance [7]. Tang and Ghorbani compared the performance of

44 HMMs with Support Vector Machine (SVM) for accent classification [8]. Kumpf and King proposed to use linear
45 discriminant analysis (LDA) for identification of three accents in Australian English [9].

46 Artificial neural networks, especially Deep Neural Networks (DNNs) and Recurrent Neural Networks (RNNs)
47 and CNNs have been widely used in state-of-the-art speech systems and Image Processing Systems [10] [11]
48 [12] [13]; however, in the area of accent identification, there are only a few studies evaluating the performance
49 of neural networks [14] [15]. Nonetheless, in a related area, language identification (LID), neural networks have
50 been investigated exhaustively [16] [17] [18]. In a recent paper [19], where they used spectrograms for accent
51 classification and speaker recognition and achieved an accuracy of 92%. Inspired by their work, we also propose
52 to use Mel-Spectrograms and MFCCs for our research

53 The rest of the paper is organized as follows: In section 2, we discuss the collection and making of dataset. In
54 section 3, we discuss proposed system and discuss in detail the features that we have used for our research. In
55 section 4, we discuss the experiments and show our results and finally in section 5, we conclude our research.

2 II.

3 Dataset a) Collection of Dataset

58 Data is very important in every machine learning and deep learning project or deep learning research. The
59 data, we required for our research was, the audio files of people speaking some sentences that we choose. These
60 sentences, to some extent captured a wide range of accent changes in the spoken Kashmiri Language. In total,
61 20 sentences were chosen for research purposes and people were recorded speaking these sentences in their native
62 accents of Kashmir language. The data was collected from 5 districts or areas of Kashmir and all these files
63 were saved with the extension of 'ogg', which in preprocessing, were changed to 'wav' format. In total, we got
64 almost 100 voice samples from each area and thus we had, 500 total voice samples of these sentences spoken by
65 different people.

4 b) Making of dataset

67 The data we had, were audio files and we decided on getting the MFCC and Mel-Spectrograms from these audio
68 files. So, our final dataset consisted of images of MFCC and Mel-Spectrograms. Since deep learning models
69 require huge amount of data, we had to augment the data to increase the size of our dataset. There are many
70 great techniques of augmenting the data, when it comes to images and audio. Since the images were of features,
71 we could not use the normal augmentation techniques like distortion, rotation and many more [reference]. A
72 special kind of augmentation known as specAugment [20], which produces augmented images on spectrograms
73 was used. This augmentation performed following operations on the Images of Mel-Spectrograms. 1) Frequency
74 masking is where certain part of the frequency is masked out, and 2) Time masking, where certain part of time
75 is masked out. Even though, we performed augmentation on Mel-Spectrogram images, the data was not enough.
76 so, we had to perform the augmentation on the audio files also. The audio files were augmented by increasing
77 the speed, pitch and amplitude of the audio files, thus giving somewhat variability in the initial dataset of audio
78 files.

79 After performing, such augmentations we had large sufficient dataset for deep learning.

5 III.

6 Proposed System a) Architecture

82 We used CNN based architecture with ReLU activation function for internal nodes and SoftMax function to
83 output the probability distribution of output classes. CNN [19] models show state of the art performance with
84 image data. Since our motive was to extract the features from the audio data and plot them as images and then
85 those images were input to the model, so we chose the model based on CNN architecture. Our model has, six
86 convolution layers and six Maxpooling layers, with 5 dense layers and a flatten layer.

87 [image of model] Accuracy and loss varied based on the feature used and learning rate of the model. We choose
88 different learning rates based on different features that were input to the model. Our models were trained on
89 learning rates between 0.001 to 0.0001

7 b) Features Used

91 Many features have been used in researches of audio processing. We decided to keep our research simple so we
92 decided on two features, MFCC and Mel Spectrograms. MFCC have been found to perform well in case audio
93 classification [21] purposes and Mel spectrograms and Spectrograms have also shown such performance in many
94 cases [19]. Different number of coefficients can be used, mostly 13 are taken. The selection of such number
95 of constants, depends on the problem in hand. We experimented on various number of coefficients and finally
96 decided on 13, 24 and 36 coefficients. These features were extracted and plotted as images and then such images
97 were input to our model.

98 Mel Spectrograms -A Mel spectrogram is a spectrogram that converts the frequencies to the Mel scale. When
99 the spectrogram from the audio file is plotted using Mel scale, we get the Mel-Spectrogram. These spectrograms
100 were plotted as images, same as the MFCCs and given input to the model.

101 All these operations of feature extraction were done using librosa library [12], which makes working with
102 audio very easy.

103 8 IV.

104 9 Experimental Setup And Results

105 Different experiments were performed on different features and different learning rates were set during the training
106 of the models.

107 The features, were stored in two ways -Images of the features were generated, and in other features were
108 extracted and a dimension was added, no image was generated and the features were stored in JSON data
109 format. Below, we show the results of various experiments:a) Experiment 1

110 This Experiment was done using images of Mel spectrograms and MFCC for the CNN having 3 color channels.
111 In this experiment, the images were generated from the audio files. Those images were saved and later were
112 loaded back into the model. The models were trained on those images and evaluated on the validation and
113 testing sets.

114 10 i. Mel Spectrograms

115 The below figure shows the metrics graphically and we can conclude from the graph that the model is showing
116 state of art results on our data. These above results were calculated on the testing data and we can conclude
117 that our model performed much better than expected and showed state of the art performance on our data.

118 11 ii. MFCCs

119 The MFCC features were extracted from the audio files and plotted as images and these Images were saved and
120 loaded at the time of model training. The following figures show the accuracies and losses with respect to the
121 epochs. Three types of constants were extracted and same model was trained on these images generated from
122 the audio files. The training was done using the training data, and validated on validated. Following table shows
123 our results. From the above table, we can see that the 24 constants performed slightly better than the others on
124 validation data.

125 Using the images as input to the model, Mel-Spectrograms and MFCC 24 constant features performed better
126 than the other features.

127 12 b) Experiment 2

128 This experiment was done using Json files of extracted features and giving to CNN having 1 color channel.

129 In this experiment, the features were extracted and were saved in JSON files. No Images were generated for
130 this data. Then the features were loaded back and the model was trained on these features. The below table
131 show the Testing accuracies and testing Losses for various features extracted from audio files. This experiment
132 was done by splitting the audio files in the chunks of 2 seconds.

133 In this experiment the audio was splitted into two second chunks a The Mel spectrograms features were
134 extracted from the splitted audio files and then saved as images and as well as JSON files. Following accuracies
135 and Losses were calculated on validation data.

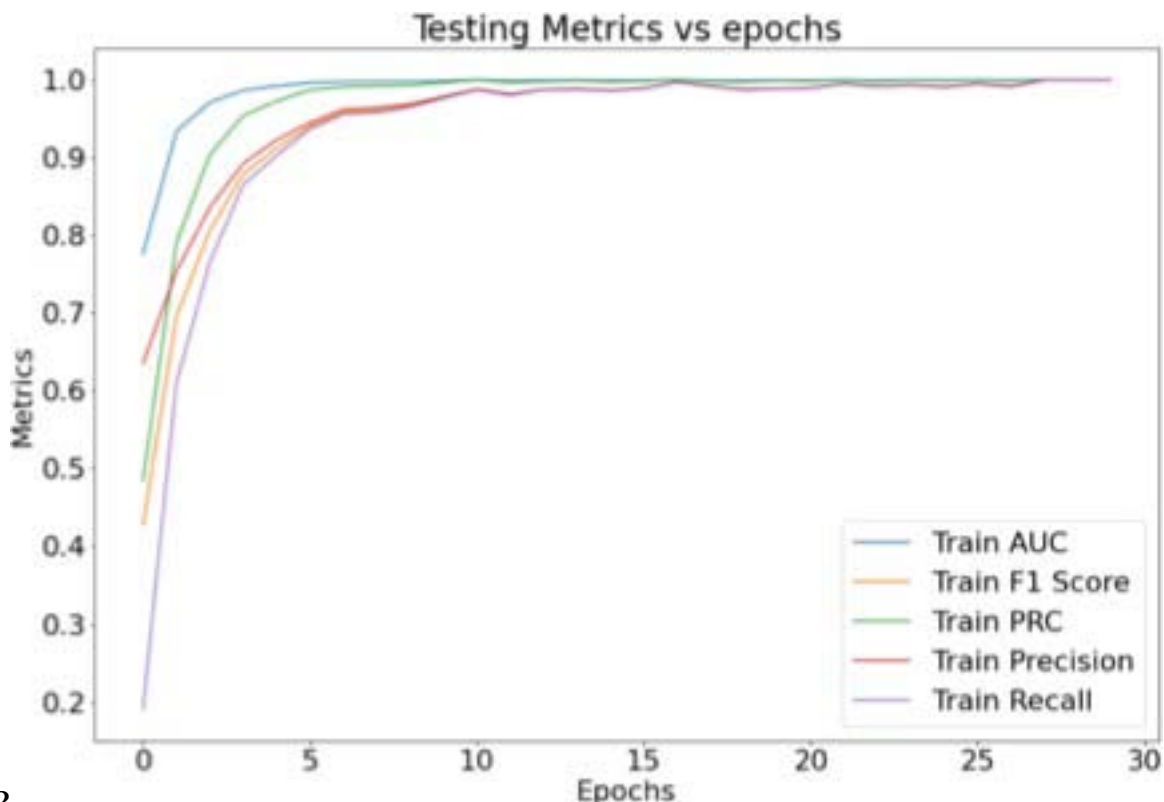
136 13 Conclusion

137 This research paper proposes a solution to the accent classification for Kashmiri Language using Convolutional
138 Neural networks. The solution is based on deep learning techniques using CNNs that adapt to the multi-
139 dimensional data. CNNs provide solution to this problem using the supervised approach where they first undergo
140 training session during which they are fed with labelled data from which they learn the relationships in the data
141 and attain the learning capability. In later stage, they are presented with unseen data of same domain and
142 are able to make remarkable inferences from this unseen data by utilizing the attained learning capability. In
143 addition to reporting the state-of-art classification results, its accuracy is also remarkable. In our research, we
144 can conclude that the models and the data that we used, the Mel Spectrograms performed better and showed
145 better performance than the MFCCs, also we saw that the images with 3 color channels performed better than
146 the features that were saved with extended dimension. Overall we can conclude, our model showed state of the
147 art performance for the Accent classification of Kashmiri Language with five output classes.

148 14 VI.

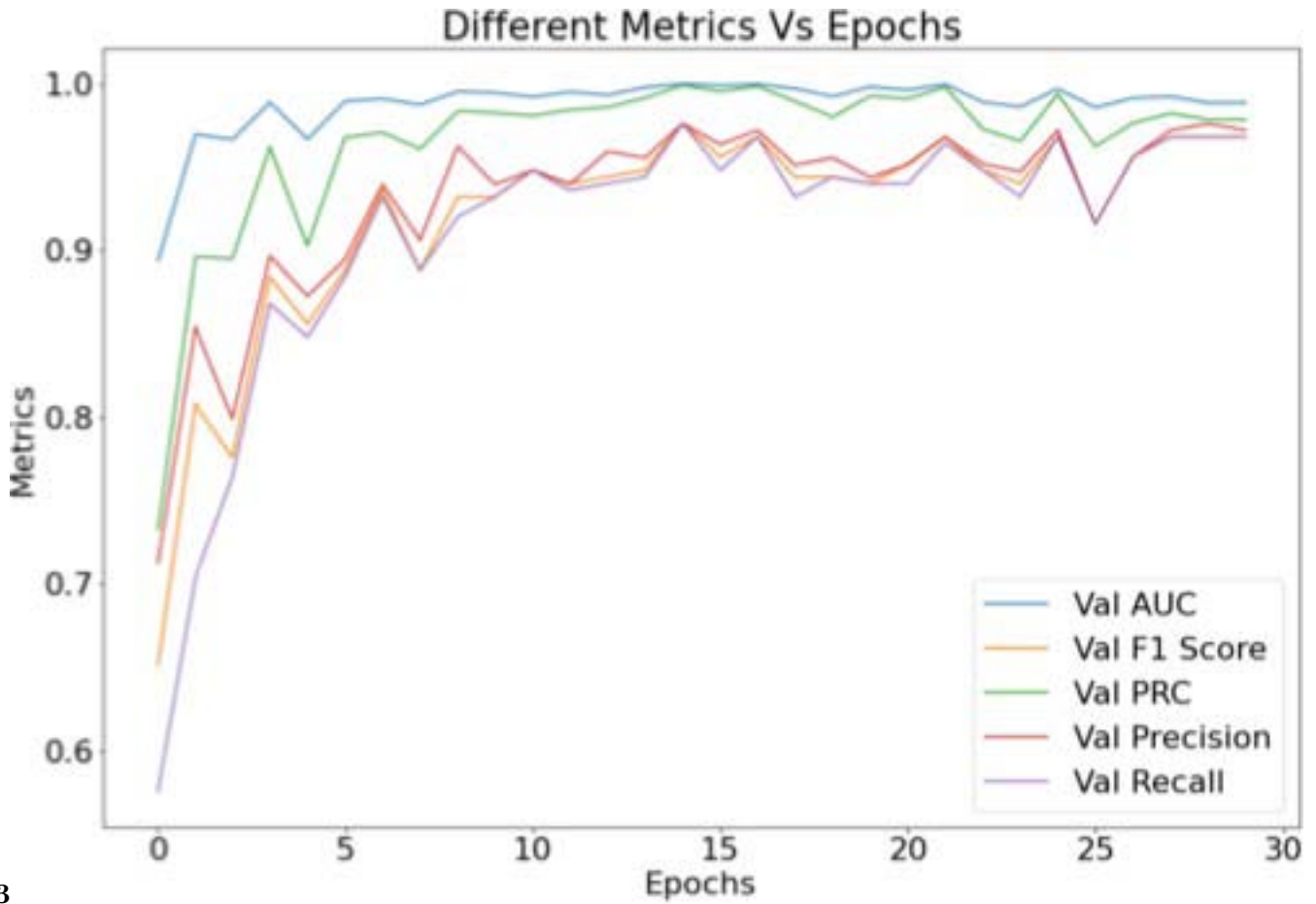
149 15 Future Improvements

150 There is a lot of improvement to be done in this field of research. Since, this is the first research in this language,
 151 as we could not find any other research related to Kashmiri language, so the area of improvement is vast. We
 152 propose following enhancements for this research-? Collection of more data for efficient model training a. The
 153 dataset can be increased in size b. The dataset can be in such a way, that it captures the maximum of the features
 154 and variations present in the language. ? The model can be made more complex and more sophisticated that
 155 would be able to handle more data and not underfit. ? Making efficient model for being able to capture most
 156 of the features of accent classification ? Improving the classification error and thus being able to classify wide
 157 range of the language. ? Making use of different architectures and techniques available for making the overall
 application most fruitful. ? The classification classes can be increased to more than 5 accents or regions.



12

Figure 1: Fig. 1 :Fig. 2 :



3

Figure 2: Fig. 3 :

1

Feature	Validation Loss	Validation Accuracy
Mel Spectrograms	0.0392	0.9848
MFCC 13	0.04	0.98
MFCC 24	0.031	0.99
MFCC 36	0.06	0.97

Figure 3: Table 1 :

2

Feature	Testing Loss	Testing Accuracy
MFCC 13	0.086	0.87
MFCC 24	0.110	0.87
MFCC 36	0.507	0.865

c) Experiment 3

Figure 4: Table 2 :

3

Type of Feature	Validation Loss	Validation Accuracy
Images	0.0331	0.98
JSON files	0.069	0.97
V.		

Figure 5: Table 3 :

-
- 159 [Hinton et al.] , G Hinton , L Deng , D Yu , G E Dahl , A .
160 [Deshpande et al. ()] ‘Accent classification in speech’. S Deshpande , S Chikkerur , V Govindaraju . *Automatic*
161 *Identification Advanced Technologies, Fourth IEEE Workshop on*, (Buffalo, NY, USA) 2005. IEEE. p. .
162 [Tang and Ghorbani ()] ‘Accent classification using support vector machine and hidden markov model’. H Tang
163 , A A Ghorbani . *Advances in Artificial Intelligence*, 2003. Springer. p. .
164 [Huang et al. ()] ‘Accent issues in large vocabulary continuous speech recognition’. C Huang , T Chen , E Chang
165 . *International Journal of Speech Technology* 2004. 7 (2-3) p. .
166 [Xu et al. ()] ‘An experimental study on speech enhancement based on deep neural networks’. Y Xu , J Du ,
167 L.-R Dai , C.-H Lee . *Signal Processing Letters, IEEE* 2014. 21 (1) p. .
168 [Chen et al. ()] ‘Automatic accent identification using gaussian mixture models’. T Chen , C Huang , E Chang ,
169 J Wang . *Automatic Speech Recognition and Understanding*, (Italy) 2001. IEEE. p. .
170 [Lopez-Moreno et al. ()] ‘Automatic language identification using deep neural networks’. J Lopez-Moreno , O
171 Gonzalez-Dominguez , D Plchot , J Martinez , P Gonzalez-Rodriguez , Moreno . *Acoustics, Speech and Signal*
172 *Processing (ICASSP), IEEE International Conference on*, (Florence, Italy) 2014. IEEE. p. .
173 [Chan et al. ()] ‘Classification of speech accents with neural networks’. M V Chan , X Feng , J A Heinen , R
174 J Niederjohn . *Neural Networks, IEEE World Congress on Computational Intelligence., IEEE International*
175 *Conference on*, 1994. IEEE. 7 p. 44834486.
176 [Montavon ()] ‘Deep learning for spoken language identification’. G Montavon . *NIPS Workshop on deep learning*
177 *for speech recognition and related applications*, (Whistler, BC, Canada) 2009. p. .
178 [Mohamed et al. ()] ‘Deep neural networks for acoustic modeling in speech recognition: The shared views of four
179 research groups’. N Mohamed , A Jaitly , V Senior , P Vanhoucke , T N Nguyen , Sainath . *Signal Processing*
180 *Magazine* 2012. 29 (6) p. . (IEEE)
181 [Biadsy et al. ()] *Dialect and accent recognition using phoneticsegmentation supervectors*, F Biadsy , J B
182 Hirschberg , D P Ellis . 2011.
183 [Zhao et al. ()] *Environmental sound classification based on feature fusion*, Huimin & Zhao , Xianglin , & Huang
184 , Wei , Lifang Liu & Yang . 2018.
185 [Kat and Fung ()] ‘Fast accent identification and accented speech recognition’. L Kat , P Fung . *Acoustics, Speech,*
186 *and Signal Processing*, (Phoenix, AZ, USA) 1999. IEEE. 1 p. . (IEEE International Conference on)
187 [Kumpf and King ()] ‘Foreign speaker accent classification using phoneme-dependent accent discrimination
188 models and comparisons with human perception benchmarks’. K Kumpf , R W King . *Proc. EuroSpeech,*
189 (EuroSpeech) 1997. 4 p. .
190 [Tanner and Tanner ()] *Forensic aspects of speech patterns: voice prints, speaker pro filing, lie and intoxication*
191 *detection*, D C Tanner , M E Tanner . 2004. Lawyers & Judges Publishing Company.
192 [Cole et al. ()] ‘Language identification with neural networks: a feasibility study’. R A Cole , J W Inouye , Y
193 K Muthusamy , M Gopalakrishnan . *Computers and Signal Processing* 1989. IEEE. p. . (IEEE Pacific Rim
194 Conference on)
195 [Mcfee ()] *librosa: Audio and Music Signal Analysis in Python*, B Mcfee . 10.1051/mateconf/201817303059. doi:
196 10/gf4wxc. MATECWebofConferences.173.03059.10.1051/mateconf/201817303059 2015. Austin,
197 Texas. p. .
198 [Jiao et al. ()] ‘Online speaking rate estimation using recurrent neural networks’. Y Jiao , M Tu , V Berisha , J
199 Liss . *Acoustics, Speech and Signal Processing, IEEE International Conference on*, (Shanghai, China) 2016.
200 IEEE.
201 [Rabiee and Setayeshi ()] ‘Persian accents identification using an adaptive neural network’. S Rabiee , Setayeshi
202 . *Second International Workshop on Education Technology and Computer Science*, (Wuhan, China) 2010.
203 IEEE. p. .
204 [Park (2019)] *SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition*, D S Park
205 . doi: 10/ghbzt4. Sep. 2019. p. . (Interspeech 2019)
206 [Zeng et al. (2019)] ‘Spectrogram based multi-task audio classification’. Y Zeng , H Mao , D Peng , Z Yi . doi:
207 10/gnkrss. *Multimed Tools Appl* Feb. 2019. 78 (3) p. .
208 [Zen and Sak ()] ‘Unidirectional long short-term memory recurrent neural network with recurrent output layer
209 for low-latency speech synthesis’. H Zen , H Sak . *Acoustics, Speech and Signal Processing (ICASSP), IEEE*
210 *International Conference on*, (Brisbane, Australia) 2015. IEEE. p. .