



## A Comprehensive Review on Phoneme Classification in ML Models

A Sai Sarath, Dr. M Sreedevi

<sup>1</sup>PG Scholar, Sri Venkateswara University, Tirupati, Andhra Pradesh, India

<sup>2</sup>Assistant Professor, Department of Computer Science, Sri Venkateswara University, Tirupati, Andhra Pradesh, India

### ABSTRACT

#### Article Info

Volume 4, Issue 10

Page Number : 131-137

Publication Issue :

July-2020

This paper gives a relative performance examination of both shallow and profound machine learning classifiers for speech recognition errands utilizing outline level phoneme classification. Phoneme recognition is as yet a principal and similarly significant introductory advance toward automatic speech recognition (ASR) frameworks. Frequently regular classifiers perform outstandingly well on domain-explicit ASR frameworks having a constrained arrangement of jargon and preparing information as opposed to profound learning draws near. It is consequently basic to assess the performance of a framework utilizing profound artificial systems regarding effectively perceiving nuclear speech units, i.e., phonemes right now customary cutting-edge machine learning classifiers. Two profound learning models - DNN and LSTM with numerous arrangement structures by changing the quantity of layers and the quantity of neurons in each layer on the OLLO speech corpora alongside with six shallow machines get the hang of ing classifiers for Filterbank acoustic features are completely considered. Moreover, features with three and ten edges transient setting are registered and contrasted and no-setting features for various models. The classifier's performance is assessed as far as accuracy, review, and F1 score for 14 consonants and 10 vowels classes for 10 speakers with 4 distinct tongues. High classification precision of 93% and 95% F1 score is gotten with DNN and LSTM organizes separately on setting subordinate features for 3-shrouded layers containing 1024 hubs each. SVM shockingly acquired even a higher classification score of 96.13% and a misclassification blunder of under 5% for consonants and 4% for vowels.

**Keywords:** Phoneme Classification, Filter-Bank, Acoustic Features, Machine Learning, SVM, DNN, LSTM, Computing Methodologies, Artificial Intelligence, Speech Recognition, Machine Learning, Feature Selection, Information Extraction, Supervised Learning, Classification.

#### Article History

Published : 20 July 2020

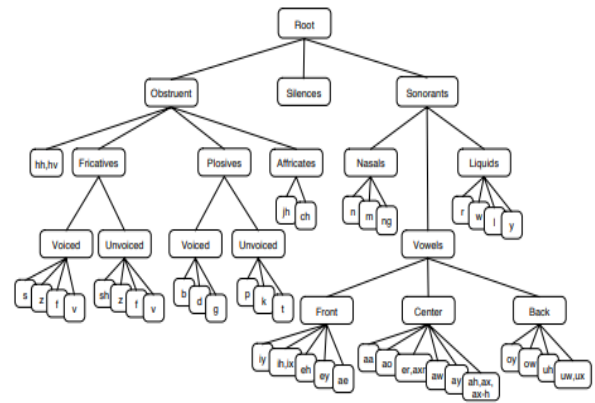
### I. INTRODUCTION

Our commitment is inspired by the way that phoneme classification at outline level can be viewed

as the front-end to the more elevated level speech recognition stage, in which the assignment of phoneme recognition by utilizing a unique programming strategy is performed, and the most

probable succession of phonemes is found. A poor front-end will altogether de-wrinkle the more significant level back-end framework. Right now, consequently, assess various procedures to find the most remarkable and appropriate model to such a classification task. Moreover, we attempt to find the best design for the profound machine learning (ML) strategies within reach, i.e., profound neural system (DNN) and long momentary memory (LSTM), by considering different basic parameters like the quantity of neurons, covered up layers, and other hyper-parameters - among others.

Phonemes classification is the undertaking of choosing what is the phonetic personality of a (commonly short) speech expression. Work in speech recognition and specifically phoneme classification commonly forces the presumption that diverse classification blunders are of a similar significance. In any case, since the arrangement of phoneme are inserted in a various leveled structure a few blunders are probably going to be more middle of the road than others. For instance, it appears to be less extreme to characterize an articulation as the phoneme/oy/(as in kid) rather than/ow/(as in pontoon), than foreseeing/w/(as in way) rather than/ow/. Besides, frequently we can't expand a high-certainty forecast for a given articulation, while as yet having the option to precisely distinguish the phonetic gathering of the expression. Right now, propose and break down a hierarchal model for classification that forces an idea of "seriousness" of forecast mistakes which is as per a pre-characterized various leveled structure. Phonetic hypothesis of spoken speech implants the arrangement of phonemes of western dialects in a phonetic progression where the phonemes comprise the leaves of the tree while wide phonetic gatherings, for example, vowels and consonants, compare to interior vertices. Such phonetic trees were depicted in [1, 2]. Persuaded by this phonetic structure we propose a progressive model (delineated in Fig. 1)



**Fig. 1.** The Phonetic Tree of American English

To remove however much information as could be expected from the edge to be arranged, we model the transient development of speech by thinking about a few edges around the present edge, as recommended in [1] and [2], where the direction in the speech was displayed by thinking about long Temporal examples (TRAP). The TRAP approaches differs from our work in the manner the feature parts are organized: In this investigation we consider a few filterbank vitality (FBE) vectors going before and following the present casing and connect them along the time hub yielding an info vector length of 280 and 840 for 3 and 10 edges individually. TRAP features, then again, are created by linking a few vitality esteems (regularly 101) at each and every basic band into one section and combining all portions after some handling (normalization). The thought is that a framework prepared by a fleeting grouping of edges is more discriminative than a model utilizing a solitary edge. Be that as it may, the subsequent feature vector has indistinguishable separating abilities from the TRAP feature regardless of the request for its parts.

## II. RELATED WORK

Machine learning strategies particularly profound neural models have been assuming an undeniably huge job in speech recognition in the most recent decade [3][4][5][6][7]. Speech recognition which has

watched a change in outlook with the development of profound learning recent years is currently generally utilized in different genuine applications, for example, subtitling video substance, sans hands interfaces in autos, and home gadgets. While humans are extraordinarily acceptable at tuning in to somebody talk and transforming speech into important words, for machines that have been a test. Analysts lately, along these lines, has put a lot of effort into assessing distinctive machine learning calculations in the field of speech recognition to improve the speech recognition capacities of a framework [8][9][10][11][12].

Supervised machine learning is a sort of machine learning calculation that utilizes a referred to dataset which is perceived as the preparation dataset to make expectations. The preparation dataset incorporates input factors (X) and reaction variables(Y). From these factors, a supervised learning calculation manufactures a model that can make forecasts of the reaction variables(Y) for another dataset (testing information) that is utilized to check the exactness of a model. A case of a supervised learning issue is foreseeing whether a client will default in paying a credit or not. The information factors here can be subtleties of the client, for example, broadcast appointment utilized, month to month pay, record as a consumer, and so forth.

Supervised learning incorporates two classes of calculations: relapse and classification calculations. There's a huge contrast between the two:

*Classification* — Classification is an issue that is utilized to foresee which class an information point is a piece of which is normally a discrete worth. From the model I gave above, anticipating whether an individual is probably going to default on a credit or not is a case of a classification issue since the classes we need to foresee are discrete: "prone to pay an advance" and "not liable to pay an advance".

*Relapse* — Regression is an issue that is utilized to anticipate persistent amount yield. A constant yield variable is a genuine worth, for example, a whole number or drifting point esteem. For instance, where classification has been utilized to decide if it will rain tomorrow, a relapse calculation will be utilized to anticipate the measure of precipitation.

Our past work [13] on the casing savvy classification of Oldenburg Logatome (OLLO) database for various talk ing rate is profoundly important and of premium not just on the grounds that this paper is the continuation of the work on the FBE feature that gave the best outcomes on the KALDI toolbox [14] for the DNN, yet in addition in light of the fact that right now current usage are performed on the Tensor ow [15] for both customary and profound learning methods yielding higher precision rates. Comparative work was completed by analysts in [16], where an experimental examination of a few ordinary ML strategies was performed on 11 parallel classification issues. The creators announced that the neural systems accomplished the best performance among different methods.

The casing shrewd classification utilizing bidirectional LSTM was explored in [17]. The outcomes in that paper show that bidirectional LSTM outflanks the standard intermittent neural system (RNN) and furthermore time windowed multi-layer perceptron (MLPs). Likewise, it was referenced that the preparation time is substantially less than different techniques. Another work on the RNN is exhibited in [18]. Right now, impact of transient setting size is considered also where the bidirectional LSTM prompts a superior precision rate. In [19] then again, bolster vector machine (SVM) was proposed as a productive model to arrange the TIMIT database in outline level phonetically.

### III. SPEECH CORPUS

Many speech corpora including TIMIT, Callfriend, Moca, NIST, Switchboard, WSJ, Voxceleb exist for speech investigation errands. The vast majority of these informational collections have been intended for explicit assignments. One such informational index called OLLO is primarily made to break down varieties in talking rates. It is a speech database that contains straightforward non-sense blends of consonants (C) and vowels (V). These mixes are called logatomes. There are 150 distinctive logatomes right now, and for every blend, the external phoneme is the equivalent.

Four distinct tongues are secured by the German speakers: no vernacular, Bavarian, East Frisian and East Phalian. The database contains logatome spoken at a normal pace, trailed by change abilities, for example, 'quick', 'slow', 'uproarious', 'delicate' and 'addressing'.

These inconstancies can be assembled into three classes:

- i. talking rate (quick, slow and normal),
- ii. (ii) speaking style (question and explanation), and
- iii. (iii) talk ing effort (boisterous, delicate and normal).

Every one of 150 logatomes have been rehashed multiple times by every speaker. A similar number of male and female speakers is utilized to record the database to cover the sexual orientation change abilities. The inspecting recurrence of the articulations is 16 kHz. OLLO has generally been utilized for examination between human speech recognition (HSR) and ASR. We primarily decided to utilize this dataset for the accompanying reasons:

- a) Evaluating distinctive changeability and their impacts on the ASR frameworks is conceivable by utilizing this database.
- b) Also, OLLO may be helpful in recognizing how tongue and highlight inuence speech recognition performance.

In the accompanying trials introduced in segment 6, 10 speakers with no lingo and normal talking rate have been picked.

### IV. MACHINE LEARNING MODELS

Right now, performance of both parametric and non-parametric machine learning classifiers is assessed on the FBE features for the speech corpus portrayed in area 3. A parametric machine learning system expect that a fixed number of parameters parameterizes the information. Basically, the measurable model of parametric procedures is speci ed by a disentangled capacity through two kinds of appropriations - (a) the class earlier likelihood, and (b) the class restrictive likelihood thickness work (back) for each measurement. The non-parametric machine learning system, then again, expect no earlier parameterized information about the basic likelihood thickness work. The classifiers, right now, exclusively on the information acquired from the preparation tests alone.

Innocent Bayes (NB) is a parametric machine learning strategy applied for classification right now, non-parametric strategies applied right now choice tree (DT) and irregular woodland (RF). DT can be viewed as one of the most well known and ground-breaking calculations in machine learning. In [21], the subject of how DTs can be utilized to improve acoustic demonstrating in speech recognition is tended to. A Support Vector Machine (SVM) can be either a parametric or non-parametric strategy. Straight SVM is a parametric classifier as it contains a fixed size of parameters spoke to by the weight coefficient while non-direct SVM, then again, is a

non-parametric system and outspread premise work part bolster vector machine, known as RBF Kernel SVM, is a common case of this family. Furthermore, two boosting procedures, Gradient boosting and Ada boosting are likewise utilized right now. These boosting systems use the troupe of classifiers creating different forecasts and majority casting a ballot among the individual classifiers.

Furthermore, a MLP and a LSTM DNNs are utilized right now. A MLP is a feed-forward artificial neural system (ANN). The artificial neurons in the system register a weighted whole of its sources of info  $x_i$ , includes an inclination  $b$ , and applies an enactment work. A basic ANN is spoken to as:

$y = f(wix_i + b)$ , where  $w$  is the gauge and  $f$  is the initiation work. Most generally utilized enactment capacities are sigmoid, which is  $(z) = 1/(1 + e^{-z})$  and rectified direct units which is  $\text{ReLU}(z) = \max(0, z)$ .

The weight and inclination terms are estimated via preparing the system on the detectable information to limit the misfortune utilizing cross-entropy or mean square mistake. In a MLP, the neurons are organized into layers. These layers are completely associated which suggests that each neuron in one layer is associated with each neuron in the adjoining layer. The information and the yield layers are unmistakable layers in the system while a system may contain various shrouded layers. Normally, a system containing more than one shrouded layer is known as a profound neural system.

LSTM is a variation of the repetitive neural system (RNN). RNN is viewed as one of the most progressive calculations that exist in the realm of profound learning. What makes LSTM one of a kind and unique contrasted with DNN is that as opposed to customary DNN that is fit for retaining long haul information, the LSTM is acceptable at keeping transient memory. LSTM is viewed as one of the most

well-known answers for the disappearing slope issue with regards to RNN. In RNN the criticism association infers that the concealed hubs add to producing the yield as well as feed their substance back onto themselves. It is the reason they have a transient memory to recollect what was their substance just beforehand. Because of its structure, LSTM has demonstrated to be effective in managing the succession of arrangement issues. Moreover, LSTM all alone is equipped for demonstrating the ow of time legitimately. A major disadvantage of LSTM is, notwithstanding, the computational cost and long preparing time. All things considered, so as to contemplate the reality of transient con-message as info, we apply a similar time-windowing as applied to other ML strategies.

## V. CONCLUSION

An examination concerning how deferent machine learning classifiers perform for outline level phoneme classification errands with and without the transient setting is completed right now. Nuclear level speech classification, a back-finish of an ASR framework, could help in the general achievement of speech recognition applications. Higher precision for the back-end sys-tem will yield many less classification blunders for the ASR utilizing customary HMM-GMM speech recognition or start to finish frameworks dependent on DNN models. The decision of a classifier in this way is basic and this selection is frequently administered by the planned application use, asset accessibility, and computational cost. This paper, in this way, attempts to make sense of the best classifier as far as classification exactness and computational expense on a sensible size database with 1.5 million FBE feature vectors for preparing and testing. Six traditional machine learning classifiers and two profound learning models with various conjurations are assessed for 24 phoneme classes containing 14 consonants and 10 vowels. 96% classification precision as far as F1 score is gotten for

SVM for  $M = 10$  amazingly as opposed to DNN and LSTM with 93% and 95% individually.

## VI. REFERENCES

- [1]. P. Schwarz, P. Matejka, L. Burget, and O. Glembek, Phoneme recognizer based on long temporal context," Speech Processing Group, Faculty of Information Technology, Brno University of Technology.[Online]. Available: <http://speech.t.vutbr.cz/en/software>, 2006.
- [2]. H. Hermansky and S. Sharma, Temporal patterns (TRAPs) in ASR of noisy speech," in IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings. ICASSP99, vol. 1, pp. 289{292, March 1999.
- [3]. A. Mohamed, G. E. Dahl, and G. Hinton, Acoustic modeling using deep belief networks," Transactions on Audio, Speech and Language Processing, vol. 20, pp. 14{22, January 2012.
- [4]. G. E. Dahl, D. Yu, L. Deng, and A. Acero, Context-dependent pre-trained deep neural networks for large-vocabulary speech recognition," IEEE Transactions on Audio, Speech, and Language Processing, vol. 20, no. 1, pp. 30{42, 2012.
- [5]. B. Kingsbury, T. N. Sainath, and H. Soltau, Scalable minimum bayes risk training of deep neural network acoustic models using distributed hessian-free optimization," in 13th Annual Conference of the International Speech Communication Association (InterSpeech 2012), pp. 10{13, ISCA, September 2012.
- [6]. D. Yu, F. Seide, and G. Li, Conversational speech transcription using context-dependent deep neural networks," in Proceedings of the 29th International Conference on Machine Learning, ICML'12, pp. 1{2, Omnipress, August 2012.
- [7]. L. Deng, J. Li, J.-T. Huang, K. Yao, D. Yu, F. Seide, M. Seltzer, G. Zweig, X. He, J. Williams, Y. Gong, and A. Acero, Recent advances in deep learning for speech research at microsoft," in 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, pp. 8604{8608, 2013. Exported from <https://app.dimensions.ai> on 2018/12/18.
- [8]. J. Padmanabhan and M. J. J. Premkumar, Machine learning in automatic speech recognition: A survey," IETE Technical Review, vol. 32, no. 4, pp. 240{251, 2015.
- [9]. I. Gavat and D. Militaru, Deep learning in acoustic modeling for automatic speech recognition and understanding-an overview," in 2015 International Conference on Speech Technology and Human-Computer Dialogue (SpeD), pp. 1{8, IEEE, October 2015.
- [10]. L. Deng and J. C. Platt, Ensemble deep learning for speech recognition," in Fifteenth Annual Conference of the International Speech Communication Association (InterSpeech 2014), pp. 1915{1919, ISCA, September 2014.
- [11]. N. Jaitly, P. Nguyen, A. Senior, and V. Vanhoucke, Application of pretrained deep neural networks to large vocabulary speech recognition," in Thirteenth Annual Conference of the International Speech Communication Association (InterSpeech 2012), ISCA, September 2012.
- [12]. J. M. Baker, L. Deng, J. Glass, S. Khudanpur, C.-H. Lee, N. Morgan, and D. O'Shaughnessy, Developments and directions in speech recognition and understanding, part 1 [dsp education]," IEEE Signal Processing Magazine, vol. 26, pp. 75{80, May 2009.
- [13]. A. S. Shahrehabaki, A. S. Imran, N. Olfati, and T. Svendsen, Acoustic feature comparison for different speaking rates," in Human-Computer Interaction. Interaction Technologies (M. Kurosu, ed.), (Cham), pp. 176{189, Springer International Publishing, June 2018.
- [14]. D. Povey, A. Ghoshal, G. Boulianne, N. Goel, M. Hannemann, Y. Qian, P. Schwarz, and G. Stemmer, The KALDI speech recognition toolkit," in IEEE 2011 Workshop on Automatic

Speech Recognition and Understanding, IEEE Signal Processing Society, December 2011.

- [15]. R. Caruana and A. Niculescu-Mizil, "An empirical comparison of supervised learning algorithms," in Proceedings of the 23rd International Conference on Machine Learning, ICML '06, (New York, NY, USA), pp. 161{168, ACM, 2006.
- [16]. A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional LSTM and other neural network architectures," Neural Networks, vol. 18, pp. 602{610, July 2005.
- [17]. M. Wollmer, B. Schuller, and G. Rigoll, "Feature frame stacking in RNN-based tandem ASR systems-learned vs. prede ned context," in Twelfth Annual Conference of the International Speech Communication Association (InterSpeech 2011), pp.1233{1236, ISCA, August 2011.
- [18]. J. Salomon, S. King, and J. Salomon, "Framewise phone classification using support vector machines," in Seventh International Conference on Spoken Language Processing, pp. 2645{2648, ISCA, September 2002.

## Authors Profile



**A Sai Sarath** received Bachelor of Business Management from Sri Rayalaseema University in the year of 2014- 2017. Pursuing Master of Computer Applications from Sri Venkateswara University, Tirupati in the year of 2017-2020. Research interest in the field of Computer Science in the area of A Comprehensive Review on Phoneme Classification in ML Models.



**Dr. Mooramreddy Sreedevi**, She is working as a Senior Assistant Professor in the Dept. of Computer Science, S.V.University, Tirupati since 2007. She obtained her Ph. D. Computer Science from S.V.University, Tirupati. She acted as a Deputy Warden for women for 4 years and also acted as a Lady Representative for 2years in SVU Teachers Association, S.V.University, Tirupati. She published 56 research papers in UGC reputed journals, Participated in 30 International Conferences and 50 National conferences. She acted as a Resource person for different universities. Her current research focuses in the areas of Network Security, Data Mining, Cloud Computing and Big data analytics.