

Speech based experiments using HTK

B117164¹

Abstract—Should add something here

I. INTRODUCTION

In recent days, we have seen a steep inclination in Speech processing, synthesis and generation research and the improvement graph had been stagnant for a decade and it has taken a deep ascent after the advent of Graphical Processing units, parallelization and vectorization. Owing to these changes we are now able to generate synthesized user independent speech by Wavenets which are basically diluted ConvNets where convolutions can be trained in parallel by different processing units. This led us to generate speech in a fraction of seconds. But, we still lack in a few areas. We need specific type of data (speech) to generate that kind of a data (speech). Modeling only imitates the data that it has been trained on. Here, in this paper, we will be dealing with a couple of hypothesis on the available data and those will either be proved or disproved by the conducted experiments.

A. ASR components:

Conventional and State of the art Automatic Speech Recognition models have similar architecture except changes in methods used in different layers. Our objective is fundamentally, given a raw audio the model has to recognize the text (or noise, if trained) spoken. The architecture is as shown in figure: (insert figure here). We will cover these levels in detail in the following sections.

B. Learning models from data:

We live in a world with so much data. All we lack is the proper structure of the data. But this can be parsed and made structured by pre-processing, data transformation techniques. So, machine or deep learning models can either discriminate the data by comparing with each other data or can learn how the data was generated. Hence the name of discriminative and generative models. As we are dealing with speech data, we are facing a non-linear data pattern which was earlier solved and learned by Gaussian Mixture Models. But GMMs have limitations of only modeling fixed state data. Then came Hidden Markov Model which can handle variable length data (vector).

C. Required data:

As far as this assignment is concerned, we need digit utterances. As usual, machine learning models can perform well if the size of training data is increased. It also fits for testing data so that certain models can be regularized or generalized to model exceptional or unseen circumstances. Having said these, we will need two important features in

the data. Firstly, large variance in the data i.e., multi-speaker with different age, gender, dialect data. Secondly, accuracy of the labeling should be higher. As the saying goes, 'We can only generate data if that was already trained in the model'. A fair 70:30 ratio between training and testing data is chosen as a rule of thumb.

D. What is hard?

So, the problem is we need data from different kinds and we can only reproduce or recognize data if they have already been fed into the model for training. There are a few circumstances where we have to consider trade-offs between using all the available data or pick the best of the data and by 'best' what's the determining factor. We will be discussing about these in the experiments tried below.

II. THEORY

A. Data Collection

Here, we are dealing with two types of speech recognition model. One, a speaker dependent and an independent model. However the concept and the theory of the models are same but the data collection and feeding. ***** The language model is a generative model that, given a sentence, emits a word sequence. There are no HMMs, Gaussians or MFCCs involved in the language model.

In the assignment, you have been given a language model for the isolated digit task. This is written by hand as a grammar and then converted (by HParse) into the equivalent finite state graph.

In HVite, this finite-state language model is compiled together with the acoustic models, to make a single finite-state recognition network. The result is a generative model that, given a sentence, emits a sequence of MFCCs.

In the section of the report on language modelling, you'll need to explain this language model, and (hopefully) also your language model for connected digits. Even in the isolated digit case, the language model is very important: it constrains the space of things that the recognition network can generate to be only sentences containing a single word.

In this course, there are no algorithms specifically associated with the language model (we are not covering how to train an n-gram language model, or how to evaluate it in isolation).

The Viterbi algorithm is used both during recognition (in HVite) and in the initial phase of training (in HInit).

The Baum-Welch algorithm is used in the final phase of training (in HRest).

The Forward algorithm is not used on its own. It is part of the Baum-Welch algorithm (which can also be called Forward-Backward). In principle, we could use the forward algorithm to perform recognition, but in practice we use an approximation to it: the Viterbi algorithm.

1) *Signal Processing*: Given the raw audio, this layer removes noise and eliminates distortion of channel (less distortion if an exclusive microphone is used and more distortion if a bad microphone that picks up environment noise is used). For the dependent model, I had to record my own voice and label it for both testing and training. HResults in htk outputted a confusion matrix of

2) *Feature extraction / engineering*: 'Like, spectrograms are easy to read consonants and vowels, waveform is not easy to discern text'. So the time based waveform of audio signal is converted into a frequency based form by applying Fast Fourier Transform (FFT). For intuition, we can think of this form as a spectrogram. Hence, the quote above. To get continuous and intricate details, we segment the data into multiple frames typically for 25 ms with 10 ms gaps. This process yield us a format called Periodogram. Then we apply Discrete Fourier Transform and Mel filter banks to get the cepstral coefficients. We only keep 12 of the total coefficients as other coefficients do not help much in recognition. As we are in frequency domain, we tried to mimic human perception and implement mel-scale conversion which is a logarithmic non-linear function that concerns / magnifies (not technically) much about the low frequency data and less about the high frequency data (like a log curve).

3) *Acoustic modeling*: Given a static or dynamic vector of coefficient features, we have to apply certain machine learning models to find the matching pattern between the dictionary of phonemes and vector mapping and give a confidence score (probability for every possibility of phoneme sequence given the vector feature sequence). Acoustic model has the knowledge of both the given acoustic properties and phonetics properties of a language. This is more of a speech based machine learning task comparing to the next layer which is Language modeling.

4) *Language modeling*: Given the features it discerns between characters by holding a repository of n-gram character model. This is more of a language processing and machine learning task. And we can just go over this as this does not affect much considering the assignment. (what?????)

5) *Word prediction*: Now, score from Language model and Acoustic model is combined and used to predict the possible word.

III. EXPERIMENTS

A. Meta Data cleaning

For every machine learning problem, understanding the data is very important. So, first, the information of speakers to decide hypothesis was present as an unstructured file format. And extracting information from it was the first and foremost task which led to the formulate the hypothesis. SO, a python script was written to parse it by universalizing the

delimiter between every word. This was later converted to a csv file so that *pandas* can be used to group every type and limitations of the provided data could be visualized by *seaborn* and *matplotlib*. These visualizations which are provided below helped by showing certain groups cannot be used for conducting reliable experiments (or impossible to conduct).

B. Speaker Data cleaning (Sanity checking)

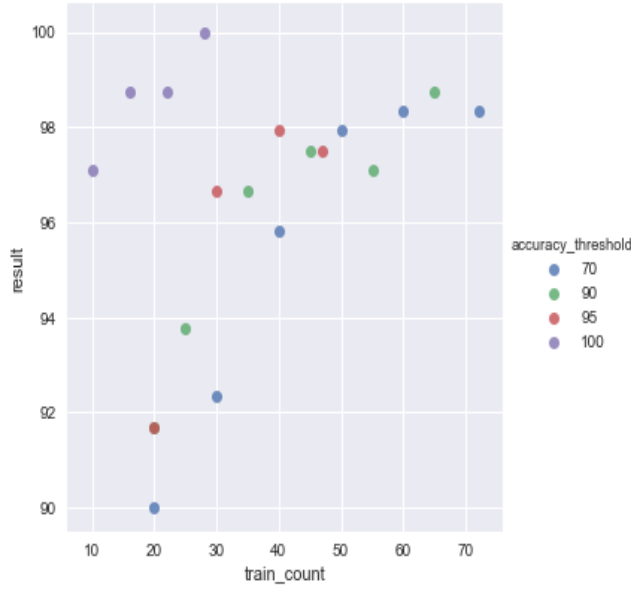
We have only considered the meta data of the speakers but the important data is the speakers audio and labeling data. After deciding the hypothesis from the sorted meta dataset and visualizations from the meta dataset, the speaker datasets are checked for extra/trivial labels like 'ten' label/data. Then for certain speakers training or testing data was not available at all. So these were detected by automating the process of running the whole experiment and the exception throwing speakers are found iteratively. Also, there we some users with few missing resources. So they were spotted out by checking the presence of 30 mfcc files per speaker and one mlf file per speaker (before concatenating all the mlf files). Few wav files for training were also missed. For checking all these exceptions, an exclusive python script was written to figure out. These inquiries were only made for the subset of the speakers who are actually involved in the experiments. Then the speakers were checked for their word error rate as an independent speech recognizer. These results were also used in the experiments.

C. Experiment I - Better data or More data? Whats better?

For production level speech recognizers, we need a huge amount of data and we cannot complete rely on the data. There are a few checks to make before using it in the model. For a more accurate model, we need more data and as it was already said we cannot trust the data completely even if the data preparation is paid for. So we do need data with good quality too i.e., properly labeled data. So, now we have a dilemma of either using all the data or only pick the high quality data. We can totally relate this scenario to the assignment.

1) *Hypothesis::* Better training data with low word error rate will yield higher accuracy and lower WER in the testing of the data.

2) *Experiment design::* To prepare the data for this experiment we need to work on a single subset where the data does not vary much in any other case like gender or accent. So it is important that we have to choose a set conforming to these properties and also the data size should be high so that we can fit multiple data size to make better observations. So to choose such a subsample, we need to pick the large samples in each category as follows.



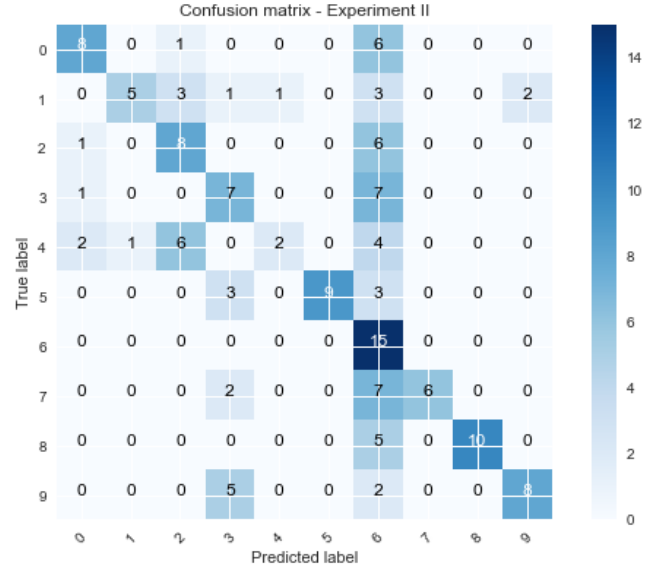
D. Experiment II - Environment noise affects the accuracy of the model

1) *Hypothesis*: Models trained on exclusive microphones (like headphones or headsets) will yield lower accuracy for inbuilt or device inclusive microphones.

2) *Experimental design*: The speakers with logitech_headsets, apple_mic, headsets were considered as the training set as they represent the exclusive microphone for better voice pick up. And the testing data were the rest of the complete data. The size of the test data was just 40 but the size of the training data was 290. The 70:30 rule was used by hindsight and the size of the training data was truncated to 90. This mixture was chose by equal distribution of accuracy rate.

3) *Result*: This concern is true with many companies as they have labeled training data with very good microphones but actual real world users may not own such microphones to use the recognizer. As expected the the accuracy for testing data was pretty bad. a 52% was scored on this experiment eventhough the training set contained many 100% accurate speaker data.

From this figure it is clear that the reason for low accuracy is due to the predicted/recognized label '6'. The reason for ambiguity can be anything. The reasons could be like as they have common pattern of phonemes or /s/ in [six] could be misinterpreted as the noise which usually high frequency fricative noise. This could correlate very much with the word [six]. This is just a theory and as this does not have anything to do with the experiments, they are left here as they have been discussed.



E. Experiment III - Digit sequences

1) *Hypothesis*: Grammars are important to define and the more sophisticated they are, better the results.

2) *Experimental design*:

F. Exploratory experiments:

Pruning variation was done by varying the factor value between 60 and 120 and a deeper knowledge of viterbi strategy of HMM algorithm was conceived. They are not added here as asked.

IV. CONCLUSIONS

APPENDIX

Appendixes should appear before the acknowledgment.

ACKNOWLEDGMENT

REFERENCES

- [1] G. O. Young, Synthetic structure of industrial plastics (Book style with paper title and editor), in *Plastics*, 2nd ed. vol. 3, J. Peters, Ed. New York: McGraw-Hill, 1964, pp. 1564.
- [2] W.-K. Chen, *Linear Networks and Systems* (Book style). Belmont, CA: Wadsworth, 1993, pp. 123135.
- [3] H. Poor, *An Introduction to Signal Detection and Estimation*. New York: Springer-Verlag, 1985, ch. 4.
- [4] B. Smith, *An approach to graphs of linear forms* (Unpublished work style), unpublished.
- [5] E. H. Miller, A note on reflector arrays (Periodical styleAccepted for publication), *IEEE Trans. Antennas Propagat.*, to be publised.
- [6] J. Wang, Fundamentals of erbium-doped fiber amplifiers arrays (Periodical styleSubmitted for publication), *IEEE J. Quantum Electron.*, submitted for publication.
- [7] C. J. Kaufman, Rocky Mountain Research Lab., Boulder, CO, private communication, May 1995.
- [8] Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, Electron spectroscopy studies on magneto-optical media and plastic substrate interfaces(Translation Journals style), *IEEE Transl. J. Magn.Jpn.*, vol. 2, Aug. 1987, pp. 740741 [Dig. 9th Annu. Conf. Magnetism Japan, 1982, p. 301].
- [9] M. Young, *The Techincal Writers Handbook*. Mill Valley, CA: University Science, 1989.
- [10] J. U. Duncombe, Infrared navigationPart I: An assessment of feasibility (Periodical style), *IEEE Trans. Electron Devices*, vol. ED-11, pp. 3439, Jan. 1959.

- [11] S. Chen, B. Mulgrew, and P. M. Grant, A clustering technique for digital communications channel equalization using radial basis function networks, *IEEE Trans. Neural Networks*, vol. 4, pp. 570-578, July 1993.
- [12] R. W. Lucky, Automatic equalization for digital communication, *Bell Syst. Tech. J.*, vol. 44, no. 4, pp. 547-588, Apr. 1965.
- [13] S. P. Bingulac, On the compatibility of adaptive controllers (Published Conference Proceedings style), in *Proc. 4th Annu. Allerton Conf. Circuits and Systems Theory*, New York, 1994, pp. 816.
- [14] G. R. Faulhaber, Design of service systems with priority reservation, in *Conf. Rec. 1995 IEEE Int. Conf. Communications*, pp. 38.
- [15] W. D. Doyle, Magnetization reversal in films with biaxial anisotropy, in *1987 Proc. INTERMAG Conf.*, pp. 2.2-12.2-6.
- [16] G. W. Juetten and L. E. Zeffanella, Radio noise currents in short sections on bundle conductors (Presented Conference Paper style), presented at the IEEE Summer power Meeting, Dallas, TX, June 22-27, 1990, Paper 90 SM 690-0 PWR.
- [17] J. G. Kreifeldt, An analysis of surface-detected EMG as an amplitude-modulated noise, presented at the 1989 Int. Conf. Medicine and Biological Engineering, Chicago, IL.
- [18] J. Williams, Narrow-band analyzer (Thesis or Dissertation style), Ph.D. dissertation, Dept. Elect. Eng., Harvard Univ., Cambridge, MA, 1993.
- [19] N. Kawasaki, Parametric study of thermal and chemical nonequilibrium nozzle flow, M.S. thesis, Dept. Electron. Eng., Osaka Univ., Osaka, Japan, 1993.
- [20] J. P. Wilkinson, Nonlinear resonant circuit devices (Patent style), U.S. Patent 3 624 12, July 16, 1990.