

Unsupervised Speaker Identification

Thivakkar Mahendran
University of Massachusetts Amherst
tmahendran@umass.edu

Varun Prasad
University of Massachusetts Amherst
varunanantha@umass.edu

1. Abstract

Filled by Varun

2. Introduction

A voice recorder makes it super convenient to take notes or preserve the minutes of a meeting. In order to make the experience better, there are tools that can be used to automatically convert speech to text. One area where this tool currently fails at is identifying the speaker. The problem that we are trying to solve is to identify the speaker. Our ultimate goal is to build a project where we are able to record a meeting and which particular person speaks at a specific time and what they speak.

It would be cumbersome and time consuming to record each speaker's voice beforehand and train a speaker recognition model on the speakers' voice. The goal is to make this tool predict the speaker without prior training on the speaker's voice. This would be an unsupervised learning project.

For example, if two speakers were talking in a meeting we would like our tool to give an output:

Speaker1: Hey, How are you?

Speaker2: I'm doing well.

Speaker2: I was able to complete the tasks that we talked about last week

Speaker1: That's great.

These two speakers' voices have not been trained before but still the program would identify the two speakers just by getting the features from their voice.

3. Background/Related Work

Filled by Varun

4. Approach

Real-world applications of such a system would require an unsupervised learning approach to the problem since it is extremely difficult to re-train the model with the voices of all the speakers in every particular situation, and also there would be no labels to train the model on. It is also difficult

to use an entirely unsupervised approach during the training phase because there would be absolutely no way to know whether the model is trained properly. Hence, evaluation of the model would become almost impossible.

Because of this, we have chosen to go with a hybrid approach where we train the model using a supervised approach over labelled data, and then use an unsupervised approach for the final classification. The first phase was done by building a fully-connected neural network and training over the 8-speaker dataset and random speaker dataset which will learn to classify the different speakers.

The second phase was implemented by cutting off the final classification layer of the neural net, and getting the embedding vector. The embedding vector at every interval, in our case 1 second, was compared to previous embeddings using a distance metric and a clustering algorithm to decide if the speaker is the same or not.

4.1. Dataset

The dataset that we initially planned to use was Vox-Celeb, which is a large scale audio-visual dataset of human speech. The dataset contains 7,000+ speakers, 1 million+ utterances, and 2,000+ hours of audio. The total size of the dataset was around 250 GB. The dataset was really huge in terms of computational complexity and also space required to store and train the model. After weeks of trying to use the dataset, we decided to build our own dataset.

We decided to create our own dataset so we could tailor the dataset to exactly suit our project. We built a pipeline to scrape audio from YouTube videos, and then split a whole audio into chunks of 1 second audio clips. We decided to create 1 second audio clips because when this model is getting used in the real world, we want to recognize the speaker instantly with as little of a delay as possible. According to the article [1] published by virtualspeech, an average person speaks about 150 words per minute, which is about 2.5 words per second. This is enough to extract useful features from the person's speech.

4.1.1 Dataset of 8 speakers

The YouTube videos that we chose to include in our dataset were speeches/monologues from celebrities. We thought this would be the best way to build a labeled dataset of audios of different people. The dataset 1 includes 7 celebrities (Obama, Hillary, Ivanka, Trump, No Speaker, Modi, Xi-Jinping, and Chadwick-Boseman) and one “no speaker” class which includes multiple background noises without anyone speaking. We included the “no speaker” class so that the model can recognize when no one is speaking. We wanted the dataset to be as diverse as possible, that is why we included speakers of both genders, different races, and also different languages. The current dataset has a size of 4.1+ hours.

	Gender	Language	Race	Length
Obama	Male	English	Black	19.5 mins
Hillary	Female	English	White	57.5 mins
Ivanka	Female	English	White	17.9 mins
Trump	Male	English	White	41.6 mins
Modi	Male	Hindi	Asian	32.4 mins
Xi-Jinping	Male	Chinese	Asian	11.18 mins
Chadwick	Male	English	Black	27.1 mins
No Speaker	N/A	N/A	N/A	39.6 mins

Table 1. Details of the 8 speaker dataset

4.1.2 Dataset of random speakers

We wanted to experiment on the quality and quantity of our dataset to see the impact on the speaker recognition model. The dataset that we created before only had 8 speakers but each speaker had multiple videos combining to more than 30 mins each. In the dataset of random speakers, we got audio files of length 2-3 minutes of **50 speakers**. In this dataset we increased the amount of speakers but limited to 1 audio file (2-3 minutes) per speaker.

The audios of the speakers were obtained from random youtube videos and voice recordings. We only chose audios which had only one speaker speaking in the whole audio file and checked if there were no background music or sound. We limited the audio file to 2-3 minutes so we could have multiple speakers. The table 2 is a snippet of the youtube videos used to create the random speaker dataset.

YouTube Title	Length
How to evaluate expressions with two variables	2:04
2020 Rock & Roll Hall of Fame Speech	2:31
A look at U.S. President-elect Joe Biden	3:32
Open Office Math	2:46
How to Get Stuff Done When You Have ADHD	4:45
How Online Math Tutoring via Skype Works	2:38

Table 2. A snippet of random speaker dataset

4.2. Supervised speaker recognition

4.2.1 Features

The first step to build a supervised speaker recognition model is to extract useful features from the audio files. These features would be the training data. We looked at multiple research papers to choose the best features for speech recognition. The features we chose, based on [2] and [3], to extract from the audio clips were:

- MFCC (Mel-Frequency Cepstral Coefficients): coefficients used to detect the envelope of audio signals. Sounds produced by humans can be accurately represented by determining the envelope of the speech signal.
- Zero-crossing rate: Number of times in a given time interval/frame that the amplitude of the speech signals passes through a value of zero. This feature is used to distinguish between periods of voiced and unvoiced sounds.
- Spectral roll off: Measure of the amount of the right-skewedness of the power spectrum. That is, the roll-off point is the frequency below which 85% of accumulated spectral magnitude is concentrated.

4.2.2 8 speaker recognition model

To build the 8 speaker recognition model we experimented with different layers and hyperparameters to get us the best result. We finally chose a 5-layer network. The first layer had the input size of 128 and the activation we chose was ReLU. The second layer has 64 as it's input with ReLU activation. The third layer has 32 as it's input with ReLU activation. The fourth layer has 16 as it's input with ReLU activation. The final layer has a size of 8 with softmax activation. The optimizer used for the model was Adam with a learning rate of $3e-4$. The loss was categorical crossentropy and the metric to analyze the model was accuracy. The total trainable parameters were 123,768.

The figure 1 represents the model built to identify the 7 speakers and 1 “no-speaker” class.

4.2.3 Random speaker recognition model

Similar to the 8 speaker recognition model, we experimented with different layers and hyperparameters to get us the best result. We finally decided to use a similar network structure that we used for the 8 speaker recognition model. Which is a 5-layer network. The first layer had the input size of 128 and the activation we chose was ReLU. The second layer has 64 as it's input with ReLU activation. The third layer has 32 as it's input with ReLU activation. The fourth layer has 16 as it's input with ReLU activation. The final

Model: "sequential_2"

Layer (type)	Output Shape	Param #
dense_5 (Dense)	(None, 11847, 128)	112768
dense_6 (Dense)	(None, 11847, 64)	8256
dense_7 (Dense)	(None, 11847, 32)	2080
dense_8 (Dense)	(None, 11847, 16)	528
dense_9 (Dense)	(None, 11847, 8)	136
Total params: 123,768		
Trainable params: 123,768		
Non-trainable params: 0		

Figure 1. 8 speaker model

layer has a size of 50 with softmax activation. The optimizer used for the model was Adam with a learning rate of $3e-4$. The loss was categorical cross entropy and the metric to analyze the model was accuracy. The total trainable parameters were 124,040. The figure 2 represents the model built to identify the 7 speakers and 1 "no-speaker" class.

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 4845, 128)	112768
dense_1 (Dense)	(None, 4845, 64)	8256
dense_2 (Dense)	(None, 4845, 32)	2080
dense_3 (Dense)	(None, 4845, 16)	528
dense_4 (Dense)	(None, 4845, 24)	408
Total params: 124,040		
Trainable params: 124,040		
Non-trainable params: 0		

Figure 2. random speaker model

We chose a similar neural network model as the 8 speaker recognition model because we wanted to see which dataset helped us recognize new speakers without prior training and didn't want to change too many variables in the experiment. We also chose the similar neural network structure because we wanted the output of the embedding vector to be the same size for 8 speaker model and random speaker model.

4.3. Using embedding vector and clustering for unsupervised speaker recognition

Once the model was trained using the training data, we saved the model locally to use the model for unsupervised speaker recognition. The input for the unsupervised speaker recognition is 1 second audio files. Similar to creating the training set for the supervised learning model, we extract the same features from each audio file.

Once the test data was ready, we passed the test data to the model. We didn't want the model to output one of the 8

speakers or the random speakers. We wanted the embedding layer of the one above the last layer. So, we popped the last layer from the model and now the output of the model is an embedding vector of size 16 for each audio file. The embedding vector looked like:

[0, 0, 0, 0, 0, 57.781567, 0, 0, 44.595768, 0, 0, 13.87579, 0, 0, 0, 39.74966]

After getting the embedding vectors of each audio file, we can use a clustering algorithm to figure out the number of clusters and which cluster each audio file belongs to. Each cluster represents a speaker. For example, after running the clustering algorithm on the embedding vectors we find out that there are 4 clusters. This means that there are in total 4 unique speakers talking.

For the clustering algorithm, we used scipy's hierarchy fclusterdata. Based on the documentation page of scipy it says that: "Clusters the original observations in the n-by-m data matrix X (n observations in m dimensions), using the euclidean distance metric to calculate distances between original observations, performs hierarchical clustering using the single linkage algorithm, and forms flat clusters using the inconsistency method with t as the cut-off threshold."

We experimented with different hyperparameters to get the best results in terms of the correct number of clusters. For the hierarchy fclusterdata algorithm we set criterion to distance, metric to euclidean, and method to centroid. We also had to choose a threshold for the clustering algorithm. We had to choose a threshold that worked best with different audio files and accurate enough to predict the right number of clusters and which cluster each audio file belonged to.

There wasn't an easy way to pick the right threshold. We ran the clustering algorithm with different thresholds and manually checked if the clustering algorithm grouped the audio files correctly. Figure 3 shows an example of checking different thresholds.

```
threshold: 60.000000, number of clusters: 4
threshold: 61.000000, number of clusters: 4
threshold: 62.000000, number of clusters: 4
threshold: 63.000000, number of clusters: 4
threshold: 64.000000, number of clusters: 4
threshold: 65.000000, number of clusters: 3
threshold: 66.000000, number of clusters: 3
threshold: 67.000000, number of clusters: 3
threshold: 68.000000, number of clusters: 3
threshold: 69.000000, number of clusters: 3
threshold: 70.000000, number of clusters: 2
threshold: 71.000000, number of clusters: 2
```

Figure 3. testing different thresholds

The threshold we picked was 65, since it was the threshold which correctly identified the number of clusters for multiple different testing datasets.

4.4. No speaker recognition model

A problem that we noticed while building the unsupervised speaker recognition model was that the clustering algorithm created a different cluster for instances when no one was speaking. Since we created a class called “no speaker”, the clustering algorithm clustered that as a speaker too.

In order to fix the problem, we had to create a no speaker recognition model, which is a binary neural network which outputs whether the audio file is a background noise without anyone speaking or if someone was speaking. In this instance there are only two classes.

MODEL DESCRIPTION - need to fill

While testing the unsupervised speaker recognition model, we first ran the 1 second audio file on the No speaker recognition model. If the model outputs “no speaker” then we just output that the audio file contains no speaker. If the model outputs “speaker” then we took that audio file and ran it through the 8 speaker recognition model or random speaker recognition model to get the embedding layer and then run the clustering algorithm to find which speaker the audio file belonged to. This fixed the issue of the clustering algorithm considering background noises as a different cluster/speaker.

4.5. Live speaker recognition

A python script has been created to get a live stream of audio data from the computer, process every 1 second audio clip by extracting features on the audio and then the features get passed to the model to get the last layer embedding vector. The new embedding vector is added to a dictionary. This dictionary contains all the previous embedding vectors. That dictionary is then used to run the clustering algorithm to find which cluster that last added embedding vector belongs to. So every 1 second, the program will output the current speaker or “no speaker” if no one is talking. This is how an unsupervised speech recognition system works in real time.

5. Experiment

Filled by Varun

6. Conclusion

Filled by Varun

References

[1] Barnard, Dom. “Average Speaking Rate and Words per Minute.” VirtualSpeech, VirtualSpeech, 20 Jan. 2018, virtualspeech.com/blog/average-speaking-rate-words-per-minute.

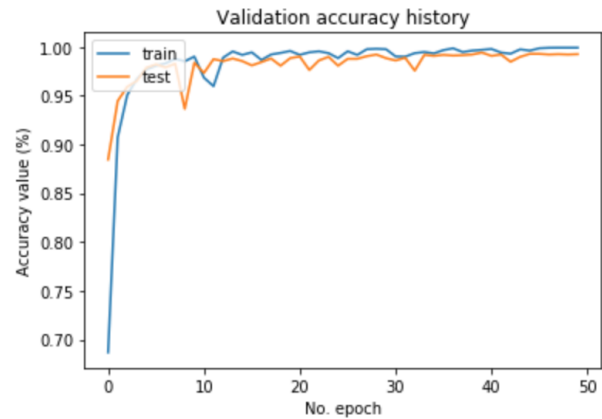


Figure 4. 8 speaker model - accuracy

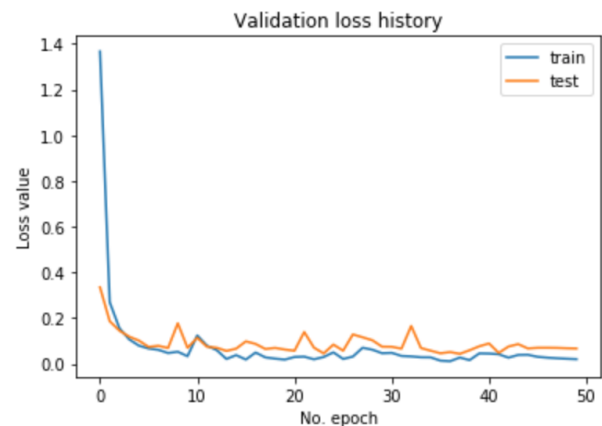


Figure 5. 8 speaker model - loss

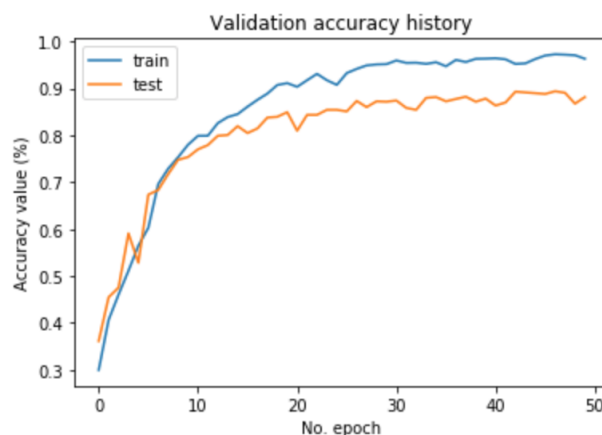


Figure 6. random speaker model - accuracy

[2] Sharma, Usha, et al. “Study of Robust Feature Extraction Techniques for Speech Recognition System.” 2015 International Conference on Futuristic Trends on Computational Analysis and Knowledge Management (ABLAZE), 2015, doi:10.1109/ablaze.2015.7154944.

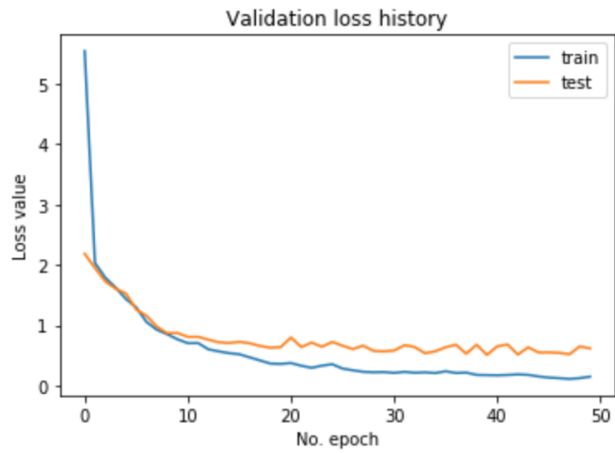


Figure 7. random speaker model - loss

- [3] “Analytical Review of Feature Extraction Technique for Automatic Speech Recognition.” *International Journal of Science and Research (IJSR)*, vol. 4, no. 11, 2015, pp. 2156–2161., doi:10.21275/v4i11.nov151681.