

Speaker Diarization and Identification Using Machine Learning

Enhong Deng, REU Student

Graduate Mentor: Abhinav Dixit, Faculty Advisors: Andreas Spanias, Visar Berisha

SenSIP Center, School of ECEE, Arizona State University

Department of Speech and Hearing Science, Arizona State University



SenSIP Algorithms and Devices REU

ABSTRACT

- Speaker diarization identifies speakers in long speech recordings.
- Form speech segments and remove undesired noise and unvoiced sections.
- Form i-vectors from features extracted from speech segments.
- Train a machine learning model on extracted features.
- Classify new speech segments according to the speaker identity.

MOTIVATION

Speaker Recognition

- Track the active speaker in a conversation with multiple speakers.

Audio Indexing

- Detect the change of speakers as a pre-processing step for automatic transcription.

Information Retrieval

- Examine contributions of speakers in speech recordings.

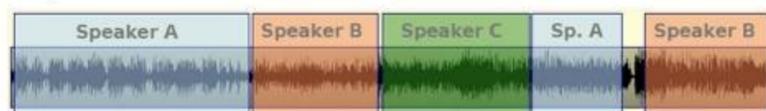
PROBLEM STATEMENT

- Perform both supervised and unsupervised speaker diarization in a telephone conversation.
- Distinguish among male and female speakers to answer the question "Who speaks when?"

Input:



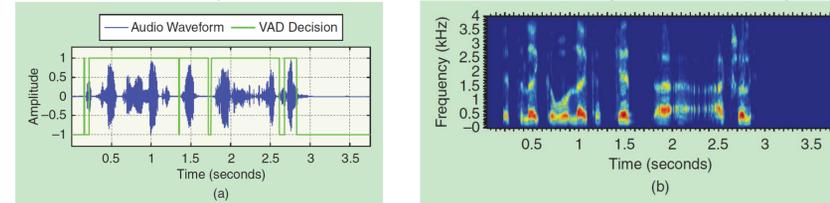
Output:



METHODS

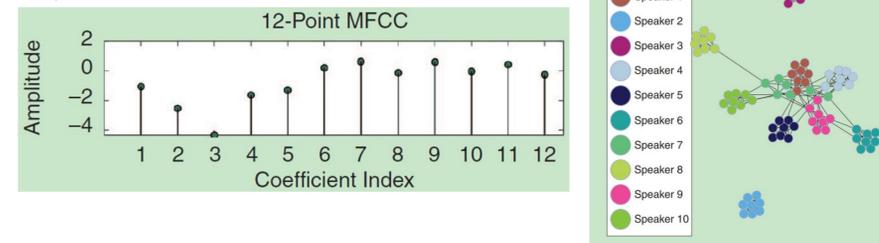
Voice-Activity Detection(VAD)

- Identifies non-speech sounds and retains only the actual speech.



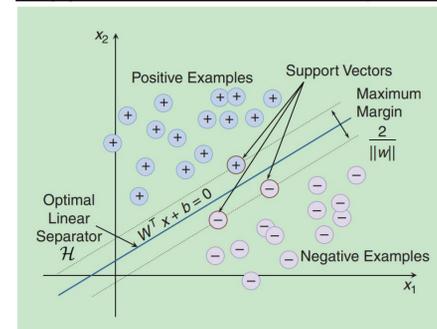
I-Vectors

- Extracting identity information using MFCCs.
- Low-dimensional i-vectors that represent the utterances from speech.



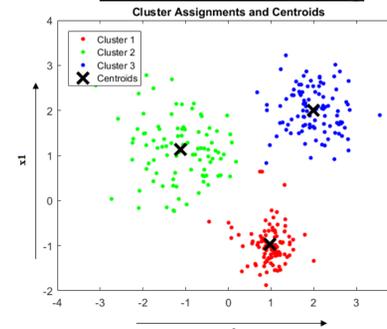
MACHINE LEARNING ALGORITHM

Support Vector Machines(SVM)



- Given labeled data, SVM can be trained to develop a model capable of distinguishing among different classes.
- The trained SVM model predicts the identity of speaker in new speech data.

K-Means Clustering



- With a known number of groups k, k number of centroids are randomly chosen.
- K-means clusters the data into k groups of clusters.

RESULTS

Supervised learning

- 97.7% accuracy in classification.
- 75% training data to generate an SVM model.
- 25% remaining data to test trained model.

Confusion Matrix

Predicted Output	Speaker 1	77 29.2%	1 0.4%	0 0.0%	98.7% 1.3%
	Speaker 2	1 0.4%	75 28.4%	1 0.4%	97.4% 2.6%
	Speaker 3	0 0.0%	3 1.1%	106 40.2%	97.2% 2.8%
	Recall	98.7% 1.3%	94.9% 5.1%	99.1% 0.9%	97.7% 2.3%
		Speaker 1	Speaker 2	Speaker 3	Precision

Speaker 1 Speaker 2 Speaker 3 Precision
True Class

Unsupervised learning

- 98.5% accuracy in clustering.
- All data is partitioned into three groups of clusters.
- Each cluster represents a speaker class.

Confusion Matrix

Predicted Output	Speaker 1	68 25.8%	1 0.4%	0 0.0%	98.6% 1.4%
	Speaker 2	1 0.4%	78 29.5%	0 0.0%	98.7% 1.3%
	Speaker 3	2 0.8%	0 0.0%	114 43.2%	98.3% 1.7%
	Recall	95.8% 4.2%	98.7% 1.3%	100% 0.0%	98.5% 1.5%
		Speaker 1	Speaker 2	Speaker 3	Precision

Speaker 1 Speaker 2 Speaker 3 Precision
True Class

- Each column of the matrix corresponds to the true class and each row corresponds to the predicted class.
- Number of correct predictions shown in green blocks, and false predictions shown in red.

REFERENCES

- J. H. L. Hansen and T. Hasan, "Speaker Recognition by Machines and Humans: A tutorial review," in *IEEE Signal Processing Magazine*, vol. 32, no. 6, pp. 74-99, Nov. 2015.
- H. Song, M. Willi, J. J. Thiagarajan, V. Berisha, and A. Spanias, "Triple Network with Attention for Speaker Diarization," in *Interspeech 2018*
- <http://multimedia.icsi.berkeley.edu/speaker-diarization/>

ACKNOWLEDGEMENT

- This project was funded in part by the National Science Foundation under Grant No. CNS 1659871 REU site: Sensors, Signal and Information Processing Devices and Algorithms.