

Speaker Recognition

SRT project of Signal Processing

Xinyu Zhou, Yuxin Wu, Tiezheng Li

Department of Computer Science and Technology
Tsinghua University

January 6, 2016

Speaker Recognition

Identification of the person who is speaking by characteristics of their voice **biometrics**.

Primary Goal:

- **Accurate** and **Efficient Short** utterance speaker recognition.

Additional Goal:

- **Scalability** over large number of speakers.
- Low **Latency** Real-Time recognition
- A working Real-Time recognition system.

Content

VAD

Voice Activity Detection shall be applied for all signals as a pre-filter. We've tried 2 different approaches:

- Energy-Based:

- Filter out the intervals with relatively low energy.
- Work perfectly for high-quality recordings.
- Sensitive to noise.

- Long-Term Spectral Divergence

- Compare long-term spectral envelope with noise spectrum.
- More robust to noise, used in our GUI.
- *Efficient voice activity detection algorithms using long-term speech information, Ramirez, Javier, 2004*

VAD

Voice Activity Detection shall be applied for all signals as a pre-filter. We've tried 2 different approaches:

- Energy-Based:

- Filter out the intervals with relatively low energy.
- Work perfectly for high-quality recordings.
- Sensitive to noise.

- Long-Term Spectral Divergence

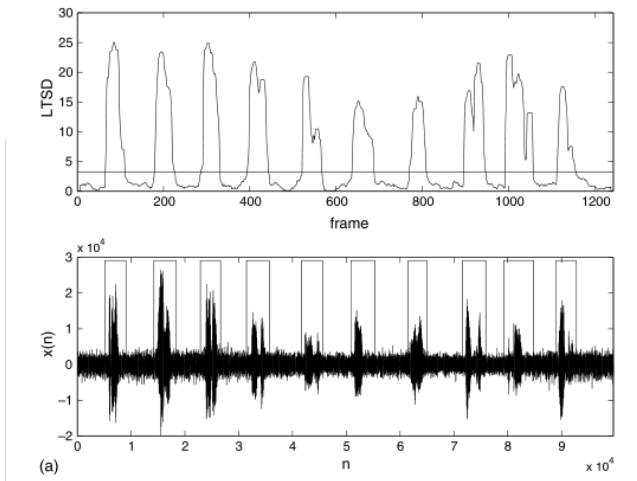
- Compare long-term spectral envelope with noise spectrum.
- More robust to noise, used in our GUI.
- *Efficient voice activity detection algorithms using long-term speech information, Ramirez, Javier, 2004*

VAD

Voice Activity Detection shall be applied for all signals as a pre-filter. We've tried 2 different approaches:

- Energy-Based:
 - Filter out the intervals with relatively low energy.
 - Work perfectly for high-quality recordings.
 - Sensitive to noise.
- Long-Term Spectral Divergence
 - Compare long-term spectral envelope with noise spectrum.
 - More robust to noise, used in our GUI.
 - *Efficient voice activity detection algorithms using long-term speech information, Ramirez, Javier, 2004*

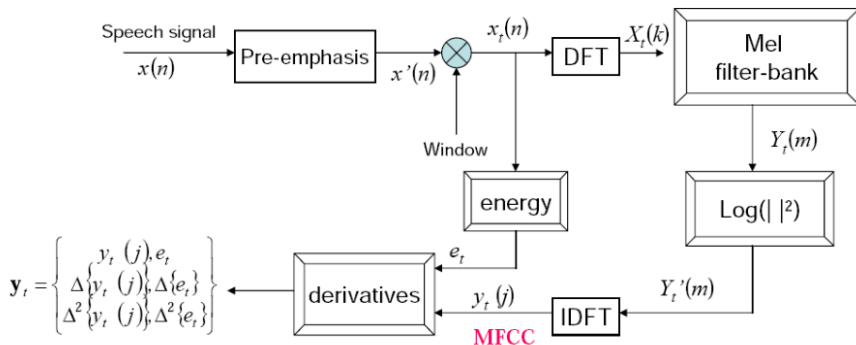
LTSD



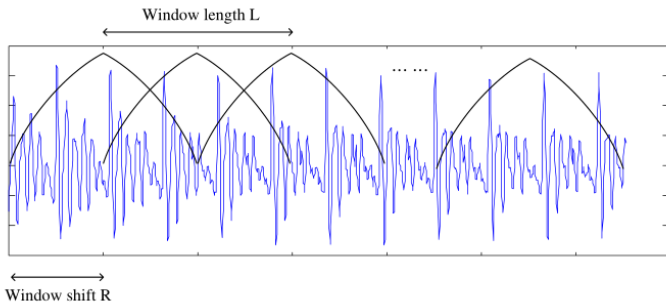
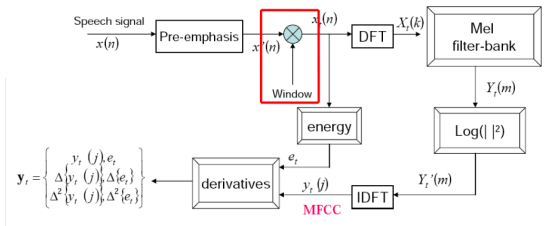
MFCC

Mel-Frequency Cepstral Coefficients

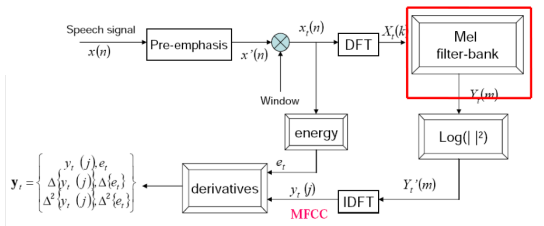
Cepstral feature which closely approximates human auditory system's response. Commonly used feature for Speech/Speaker Recognition.



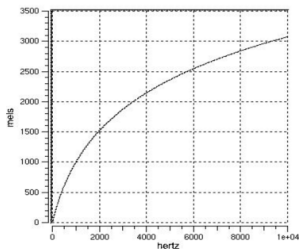
Windowing



Mel-Scale



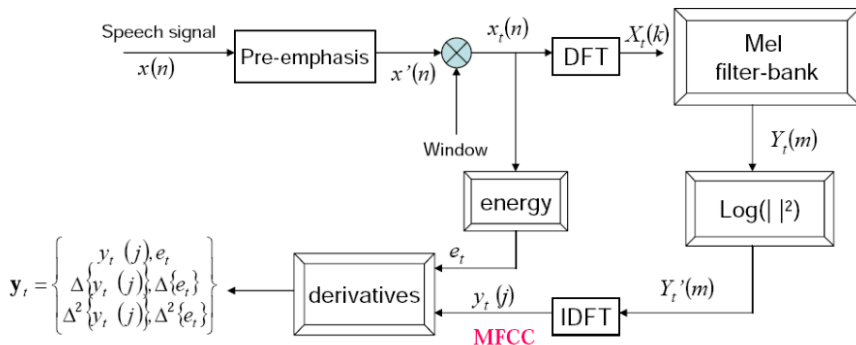
$$\text{Mel}(f) = 2595 \log_{10} \left(1 + \frac{f}{700} \right)$$



MFCC

Mel-Frequency Cepstral Coefficients

Cepstral feature which closely approximates human auditory system's response. Commonly used feature for Speech/Speaker Recognition.



LPC

Linear Predictive Coding/Coefficients

Assumption

In a short period, the n th signal is a linear combination of previous p

signals:
$$\hat{x}(n) = \sum_{i=1}^p a_i x(n-i)$$

Minimize squared error $E[\hat{x}(n) - x(n)]$ using Levinson-Durbin algorithm.

Use a_1, \dots, a_p as features.

LPC

Linear Predictive Coding/Coefficients

Assumption

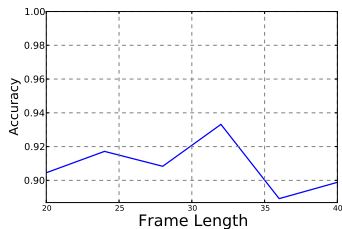
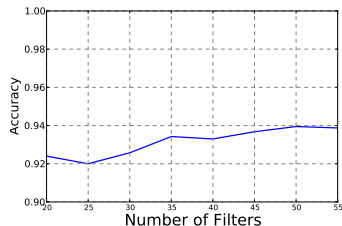
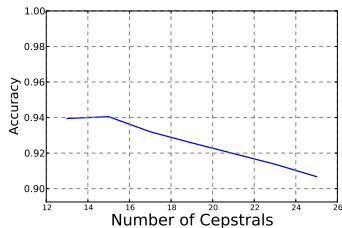
In a short period, the n th signal is a linear combination of previous p

signals:
$$\hat{x}(n) = \sum_{i=1}^p a_i x(n-i)$$

Minimize squared error $E[\hat{x}(n) - x(n)]$ using Levinson-Durbin algorithm.

Use a_1, \dots, a_p as features.

MFCC Params



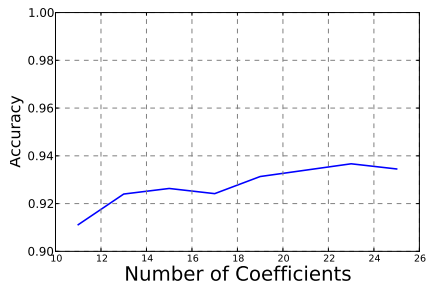
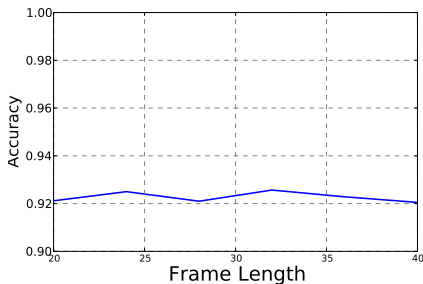
Best parameters in our cases:

Number of cepstrals: 15

Number of filters: 55

Frame length: 32ms

LPC Params



Best parameter in our cases:

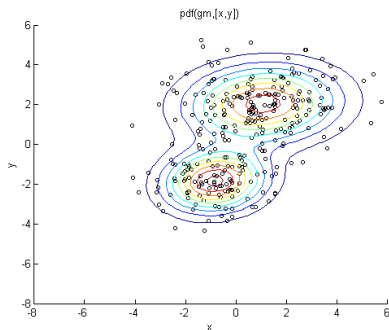
Number of coefficients: 23

Frame length: 32ms

GMM

Gaussian Mixture Model is commonly used to model human's acoustic feature.

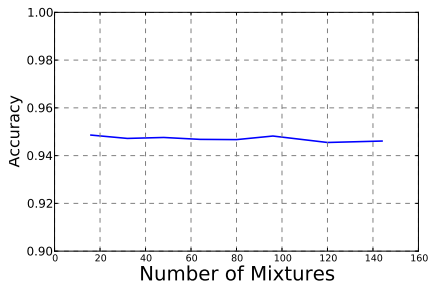
$$p(\theta) = \sum_{i=1}^K w_i \mathcal{N}(\mu_i, \Sigma_i)$$



GMM

Gaussian Mixture Model is commonly used to model human's acoustic feature.

$$p(\theta) = \sum_{i=1}^K w_i \mathcal{N}(\mu_i, \Sigma_i)$$



We use $K = 32$

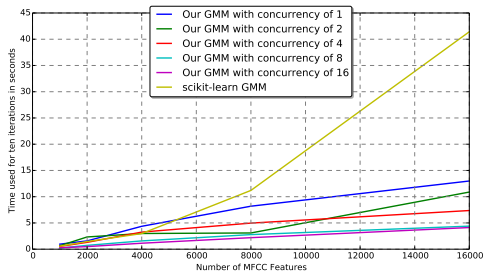
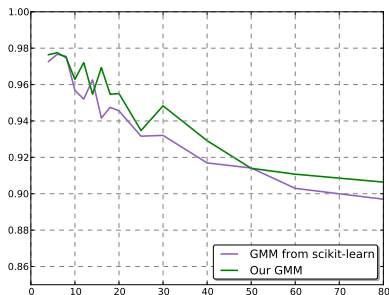
Optimized GMM

- Basic GMM training: random initialize, estimate parameters with EM.
- Improvement: initialize with a parallel KMeansII.
- Improvement: parallel training implementation in C++.
- Compared to GMM from scikit-learn:

Arthur, David, Sergei, 2007, k-means++: The advantages of careful seeding.
Bahmani, et. al, 2012, Scalable K-means++

Optimized GMM

- Basic GMM training: random initialize, estimate parameters with EM.
- Improvement: initialize with a parallel KMeansII.
- Improvement: parallel training implementation in C++.
- Compared to GMM from scikit-learn:



Arthur, David, Sergei, 2007, *k-means++: The advantages of careful seeding.*
Bahmani, et. al, 2012, *Scalable K-means++*

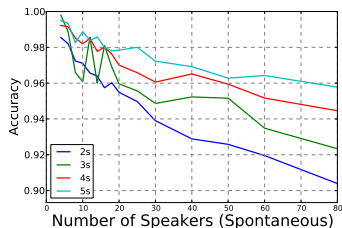
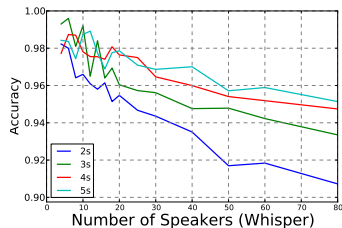
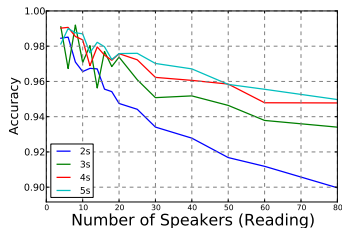
UBM

Universal Background Model is a GMM trained on giant datasets.
UBM can be used to:

- Describe general acoustic feature of human.
- Reject the decision of GMM.
- Train adaptive GMM.

*Reynolds, Douglas, et al, 2000,
Speaker verification using adapted Gaussian mixture models*

GMM Results



Train duration: 20s
Random selected test utterance: 50
Each value in the graph is an average
of 20 independent experiments.

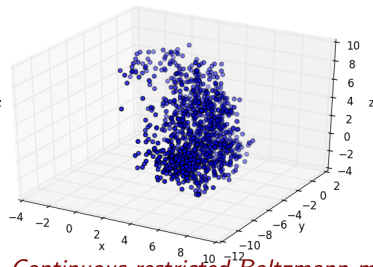
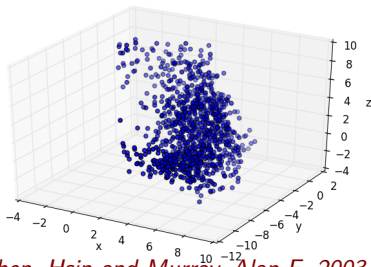
CRBM

- **Restricted Boltzmann Machine** is a generative stochastic two-layer neural network.
- **Continuous RBM** extends RBM to real-valued inputs.
- RBM has the ability to reconstruct a layer similar to input layer. The difference between the two layers can be used to measure the fitness of an input to the model.
- Therefore, RBM can be a substitution to GMM.

Chen, Hsin and Murray, Alan F, 2003, Continuous restricted Boltzmann machine with an implementable training algorithm

CRBM

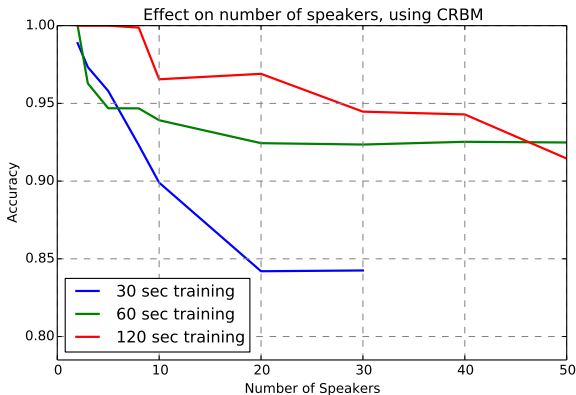
- **Restricted Boltzmann Machine** is a generative stochastic two-layer neural network.
- **Continuous RBM** extends RBM to real-valued inputs.
- RBM has the ability to reconstruct a layer similar to input layer. The difference between the two layers can be used to measure the fitness of an input to the model.
- Therefore, RBM can be a substitution to GMM.



Chen, Hsin and Murray, Alan F, 2003, Continuous restricted Boltzmann machine with an implementable training algorithm

RBM Results

Results of CRBM, tested with 5 secs of utterance.



GUI Demo

Conclusion

- We implemented a faster GMM, also with better performance.
- Accuracy is kept even under short training and testing utterance.
- Our system is highly accurate, can almost response in real-time.
- 97% accuracy for 20~30 speakers, 95% for 70~80 speakers.

Thanks!