



Speech Recording based Language Recognition

Leopold Cambier, Cindy Orozco Bohorquez, Matan Leibovich
lcambier@stanford.edu, orozcocc@stanford.edu, matanle@stanford.edu
 Institute of Computational and Mathematical Engineering



Overview

- **Motivation:** Construct a real time language classifier for communication purposes.
- **Method:**
 - Construct ML estimator based on Gaussian Mixture Model (GMM) density estimations.
 - Feature vectors based on Shifted Delta Cepstral Coefficients (SDC)
- **Results:**
 - Classification error decreases with model complexity to a certain limit.
 - Optimal number of gaussians varies significantly across languages.
 - There is clustering pattern in the classification.

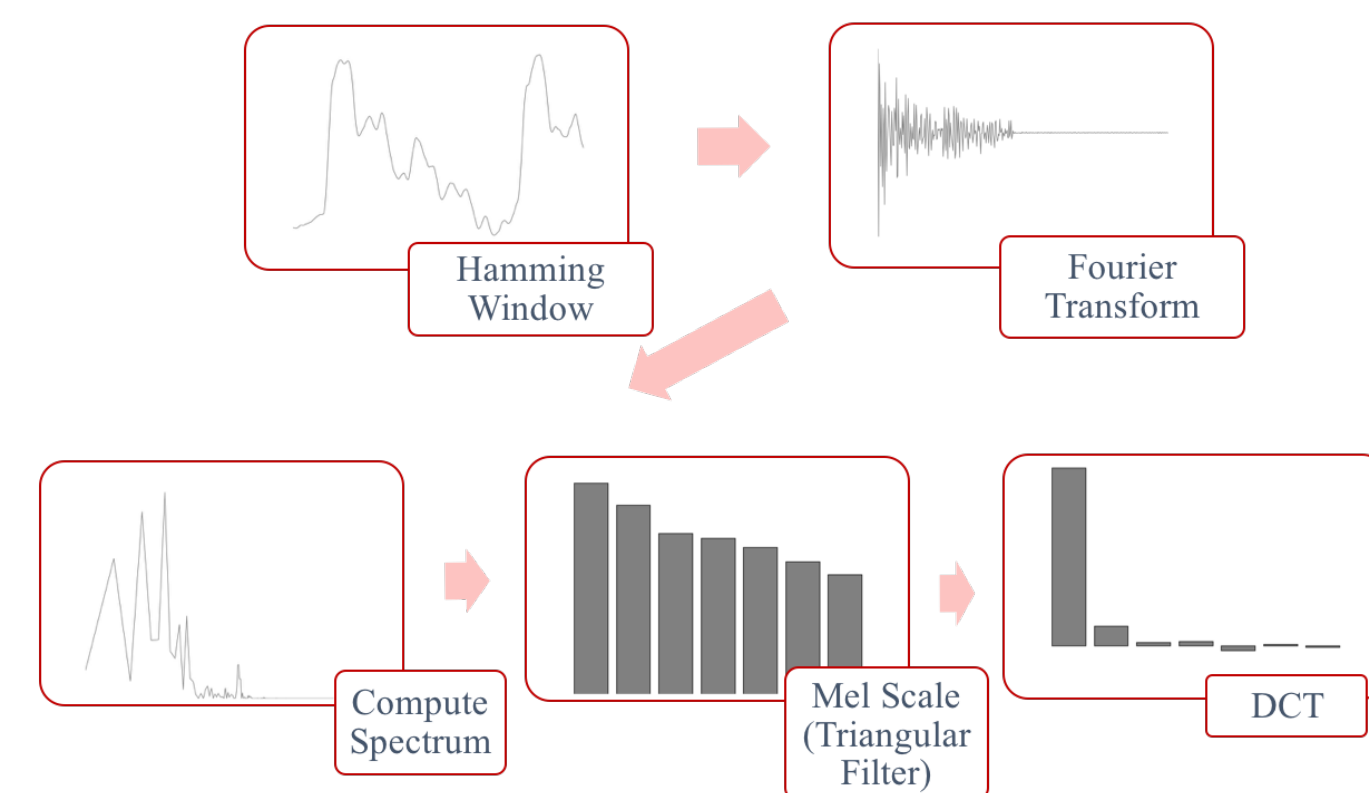
Data

Source: <https://community.topcoder.com/longcontest/?module=ViewProblemStatement&rd=16555&pm=13978>.

- \mathcal{L} : 176 different languages (some very exotic!)
- 376 10 sec. samples for each language
- Each sample divided into 210 ms long segments

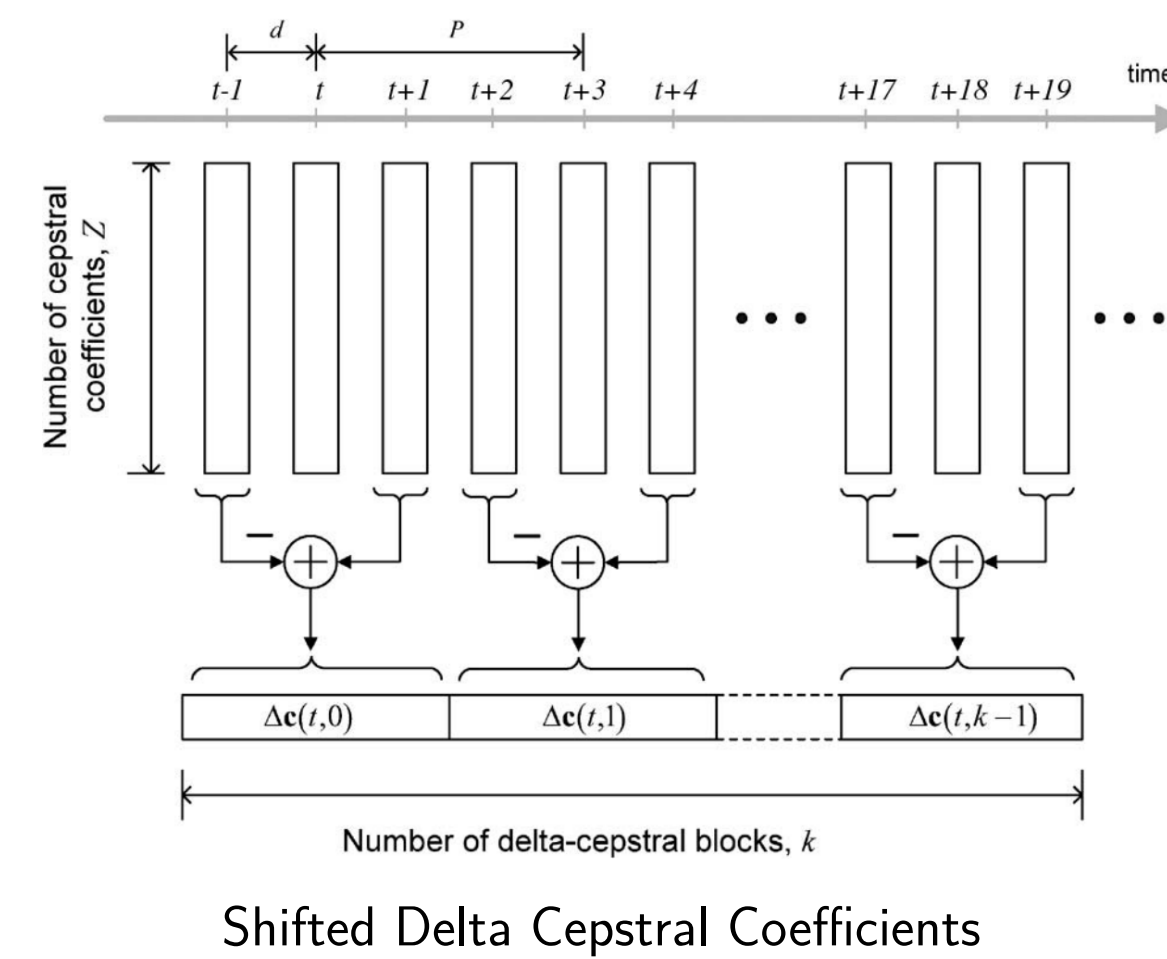
Feature Extraction: Shifted Delta Cepstral (SDC)

- 1 **Hamming Window:** $x_{\text{Ham}}(t) = (x^* \text{Ham})(t)$
- 2 **Fourier Transform:** $\hat{x}(f) = \mathcal{F}(x_{\text{Ham}}(t))$
- 3 **Transition to Mel Scale:** $\hat{x}_{\text{Mel}}(f) = \hat{x}(\mathcal{M}(f))$



Construction of Cepstral Coefficients

- 4 **DCT of logarithm:** $x_{\text{Cep}}[i] = \text{DCT}(\log \hat{x}_{\text{Mel}}(f))[i]$
- 5 **SDC:** $x_{\text{SDC}}^k[i] = x_{\text{Cep}}^{k+1}[i] - x_{\text{Cep}}^{k-1}[i]$



- 7-vectors cepstral coefficient for each 10 ms segment
- Use 3 adjacent cepstral coefficients to construct a 7-vector SDC coefficient
- For each 210ms sample \Rightarrow 49 SDC coefficients
- $46 \times 376 = 17296$ data samples $x_{\ell}^{(i)}$

Gaussian Mixture Model

Density Estimation

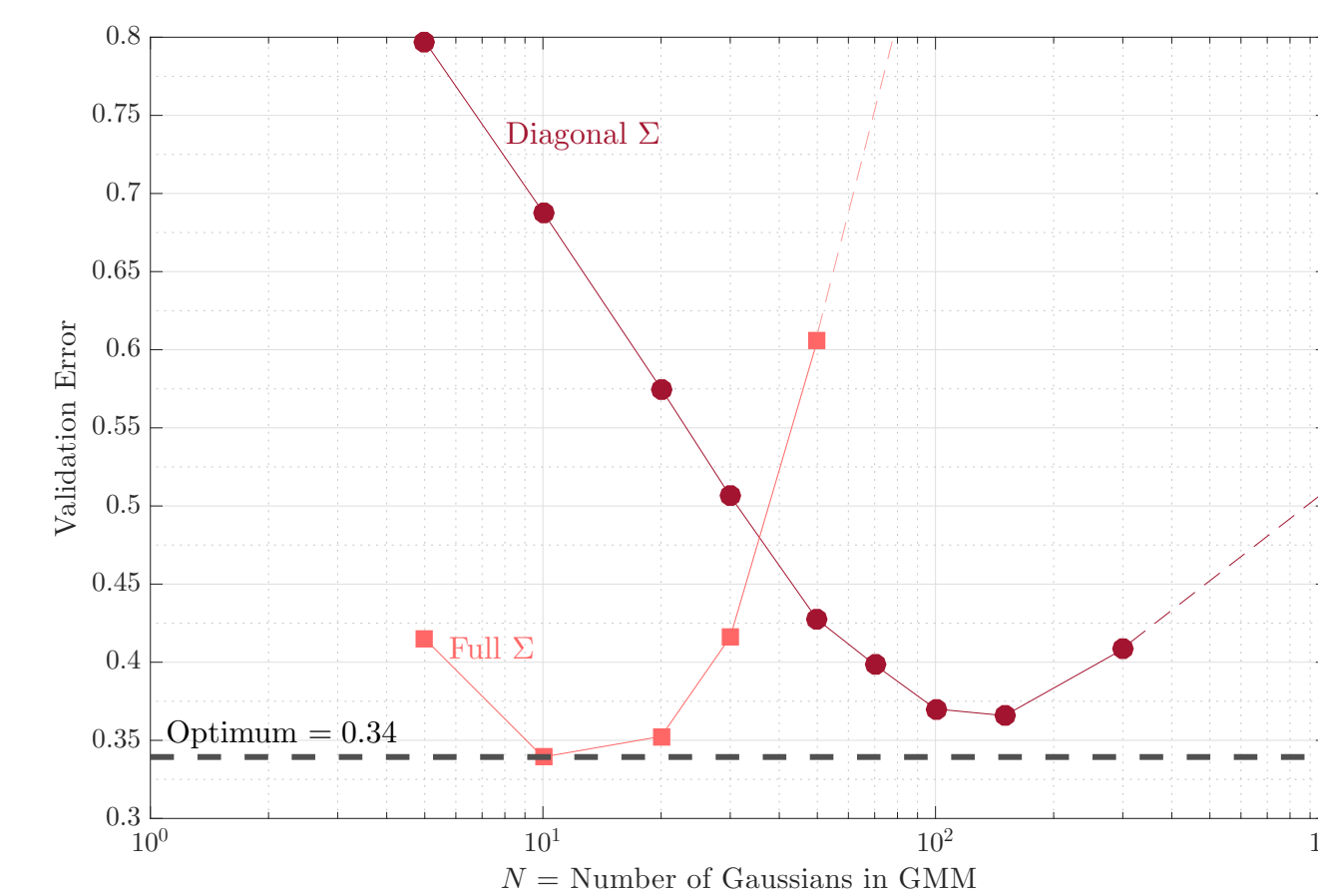
- Use $x_{\ell}^{(i)}$ to estimate mixture of N Gaussian densities
- $$p_{\ell}(x) = \sum_{k=1}^N w_{\ell}^k P(x; \mu_{\ell}^k, \Sigma_{\ell}^k), \quad x \in \mathbb{R}^{49}$$
- $$P(x; \mu_{\ell}^k, \Sigma_{\ell}^k) = \frac{1}{(2\pi^{49/2} |\Sigma_{\ell}^k|)^{49/2}} e^{-\frac{1}{2}(x - \mu_{\ell}^k)^T [\Sigma_{\ell}^k]^{-1} (x - \mu_{\ell}^k)}$$
- $$\mu_{\ell}^k \in \mathbb{R}^{49}, \Sigma_{\ell}^k \in \mathbb{R}^{49 \times 49}$$
- Σ_{ℓ}^k - either full or diag.
 - Training using Python's `sklearn` toolkit.

Classification

- Given 10s sample, construct
- $$x^i, \quad i = 1, \dots, 46$$
- SDC vectors for each 210 ms segment.
- MLE estimator
- $$\hat{\ell} = \arg \max_{\ell \in \mathcal{L}} \prod_{i=1}^{46} p_{\ell}(x_i)$$

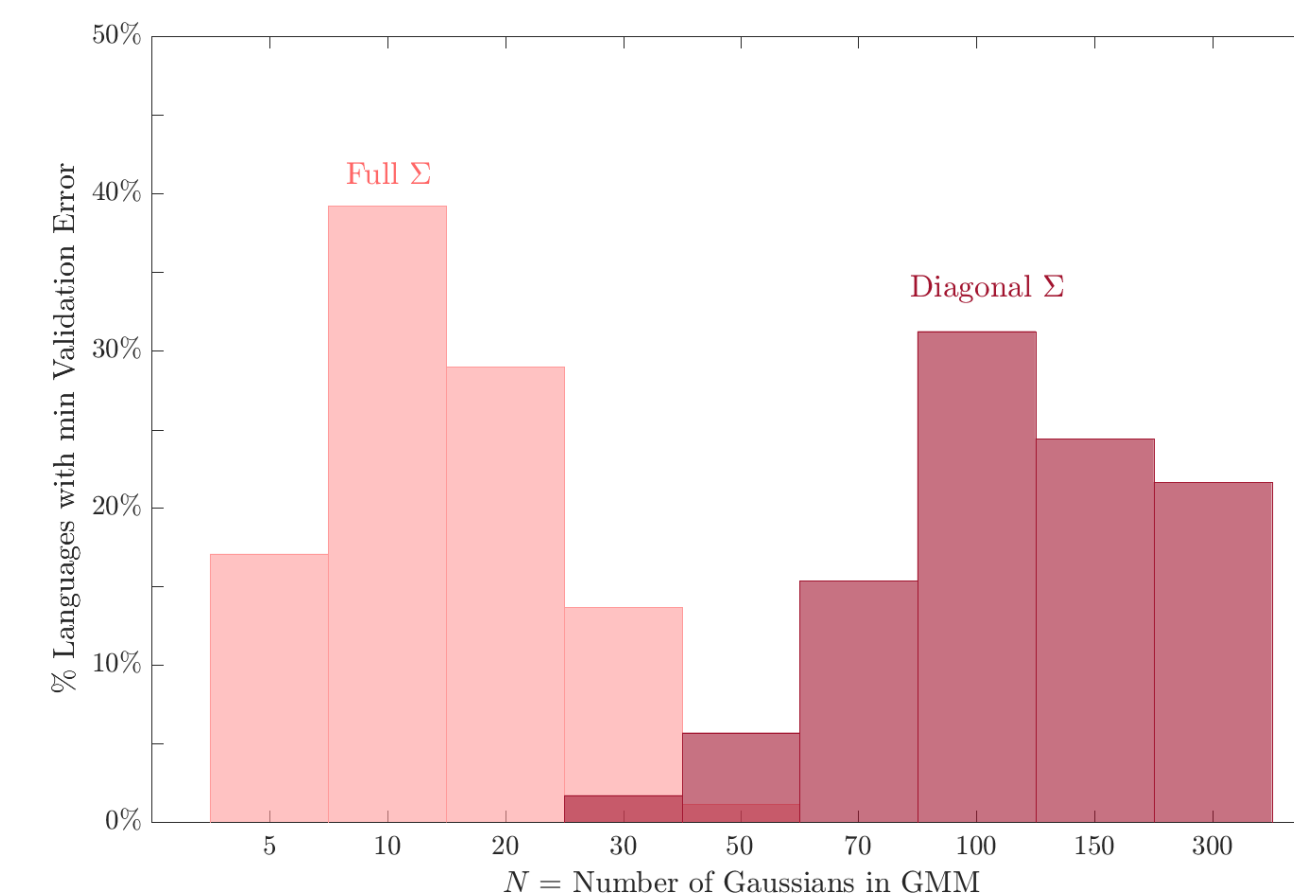
Results

- Model trained with full and diagonal covariances
- Optimal for full covariance at $N = 10$



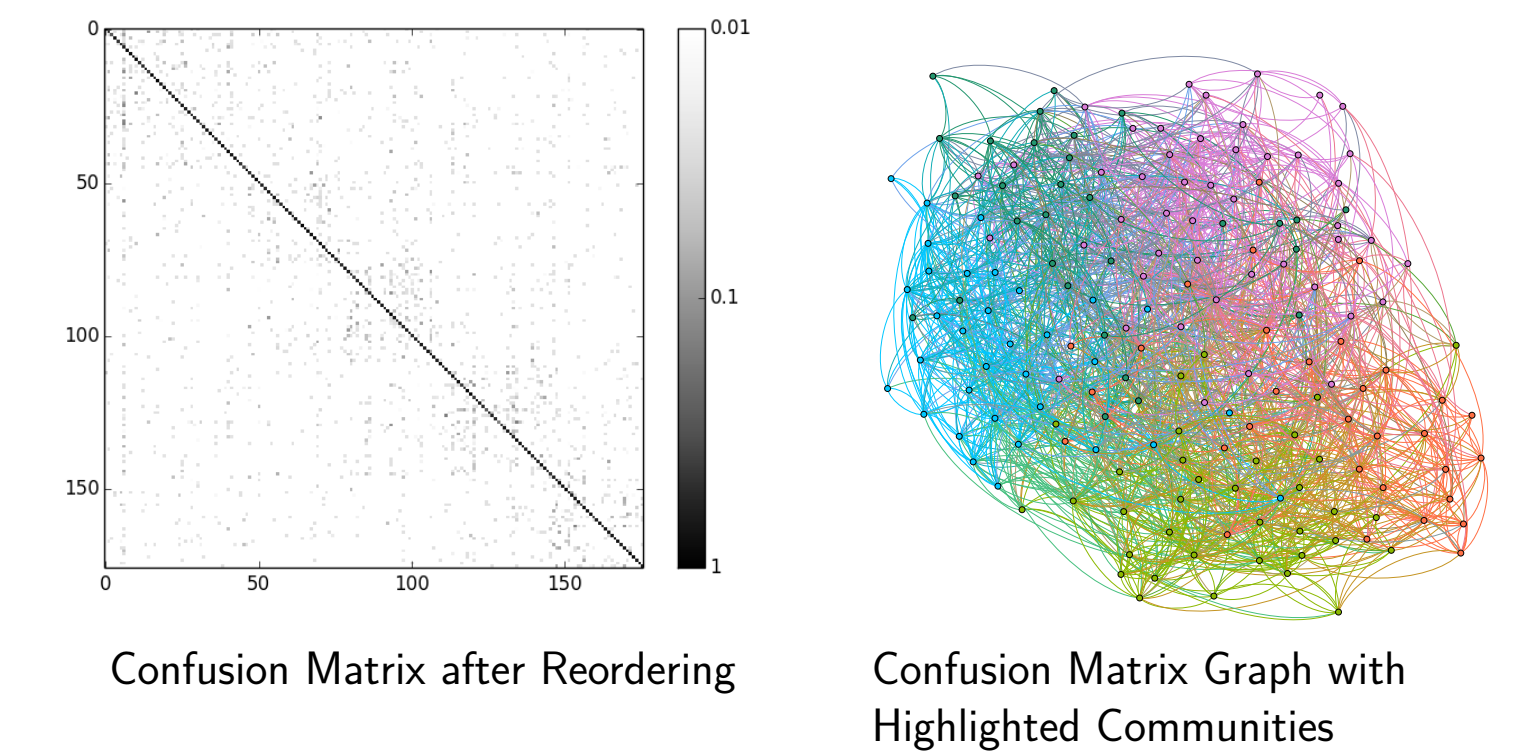
Validation Error as Function of N

- Optimal N actually varies significantly from one language to another



Distribution of Optimal N over Languages

- Error is non uniformly distributed but tends to cluster
- Analysis of confusion matrix graph shows communities



Discussion

- Test error of 0.34 with $N = 10$ and full covariances, on a dataset with 176 (balanced) classes
- Increasing N (to some extent) decreases the validation error, which shows some consistency between the model and the data
- Relatively high variability in the accuracy depending on the language

Future Work

- Allow variability in N per language for model fitting
- Introduce a measure for the quality of the GMM fit
- Train the algorithm on a larger dataset
- Predict how well would the classifier fit the sample (confidence intervals)

Bibliography

- [1] Li, Haizhou, Bin Ma, and Kong Aik Lee. "Spoken language recognition: from fundamentals to practice." Proceedings of the IEEE 101.5 (2013): 1136-1159.
- [2] Torres-Carrasquillo, Pedro A., et al. "Approaches to language identification using Gaussian mixture models and shifted delta cepstral features." INTERSPEECH. 2002.