

D E P A R T M E N T O F ELECTRICAL & COMPUTER ENGINEERING

Speech Communication Lab University of Maryland SRE-10

Daniel Garcia-Romero Xinhui Zhou Carol Espy-Wilson





Overview

- Single system submission
- Core-Core condition
- Based on Signal Coding using Overcomplete Dictionaries [1]
- Multiple encoding and scoring types



Paradigm



Mathematical description

- Linear-Gaussian observation model:
 - -Likelihood: $\eta = \Phi\beta + \epsilon$, with $\epsilon \sim \mathcal{N}(0, \Sigma N^{-1})$
 - Prior: $\beta \sim \mathcal{N}(0, I)$
 - Posterior: $\beta | \eta \sim \mathcal{N}(\mu_{\beta}, C_{\beta})$
 - with $\mu_{\beta} = C_{\beta} \Phi^{T} \Sigma^{-1} N$ and $C_{\beta} = (I + \Phi^{T} \Sigma^{-1} N \Phi)^{-1}$ – Mode of posterior: μ_{β} because Gaussian

$$\boldsymbol{\beta}_{\text{MAP}} = \arg\min_{\boldsymbol{\beta}} \frac{1}{2} \left\| \mathbf{W}^{\frac{1}{2}} (\boldsymbol{\eta} - \boldsymbol{\Phi} \boldsymbol{\beta}) \right\|_{2}^{2} + \frac{1}{2} \|\boldsymbol{\beta}\|_{2}^{2}$$



Signal Coding (SC)

 $\widehat{\mathbf{\eta}} = \mathbf{\Phi} \mathbf{\beta}_{MAP}$



Dictionary Learning (DL)

• Given R utterances (η_r, W_r) :

$$\min_{\boldsymbol{\Phi},\{\boldsymbol{\beta}_r\}}\sum_{r=1}^R \left\| \boldsymbol{W}_r^{\frac{1}{2}}(\boldsymbol{\eta}_r - \boldsymbol{\Phi}\boldsymbol{\beta}_r) \right\|_2^2 + \|\boldsymbol{\beta}_r\|_2^2.$$

- Block-coordinate descent:
 - Alternating optimization
 - Signal coding \implies Dictionary Update (DU)
 - SC is performed keeping Φ fixed
 - DU is performed by ML keeping the $\{\beta_r\}$ fixed



Scoring

• Cosine similarity with multiple metrics:

$$score = \frac{\langle \hat{\eta}_{mod}, \hat{\eta}_{test} \rangle_{W_{\#}}}{\langle \hat{\eta}_{mod}, \hat{\eta}_{mod} \rangle_{W_{\#}}^{1/2} \langle \hat{\eta}_{test}, \hat{\eta}_{test} \rangle_{W_{\#}}^{1/2}}$$
with $W_{\#} = \Sigma^{-1} N_{test}$ or $W_{\#} = \Sigma^{-1} N_{UBM}$
Linear scoring [2]:
$$score = \frac{1}{N_{TOT}} \langle \hat{\eta}_{mod}, \eta_{test} - U x_{test} \rangle_{W_{test}}$$



System configuration

- Feature extraction*:
 - 38 MFCCs (c1-c19 + Delta) every 10ms with 20ms window.
 RASTA and standardized.
 - VAD based on ASR + Energy heuristics.
- UBM*: 2048 mixtures from SRE-04,05 and 06
- Dictionary based on SRE-04, 05 and 06 and SRE08 follow-up data
 - Independent training of U and V matrices. D fixed with rel=16
 - Submitted $\mathbf{\Phi} = [\mathbf{U}_{mic}\mathbf{U}_{tel}\mathbf{D}]$, 50, 100, and 2048*38 dim
 - Late $\mathbf{\Phi} = [\mathbf{U}_{mic}\mathbf{U}_{tel}\mathbf{V}\mathbf{D}]$ with 50, 100, 300 and 2048*38 dim



Results (primary submission)



Analysis Results (ZT-norm)



Great improvement by ZT-norm: mode than 50% in some cases



Analysis Results (Scoring)



Cosine similarity slightly better but not significant



Conclusions

- Established a connection between JFA and signal coding in overcomplete dictionaries
- Mixed results between cosine similarity and linear scoring
- ZT-norm is essential
- High vocal effort quite detrimental to performance
- Set a baseline for comparing with future research:
 - L1-regularized regression
 - Discriminative dictionary learning
- Thanks to MIT-LL for providing binaries for feature extraction and UBM training.



References

- [1] D. Garcia-Romero and C. Espy-Wilson, "Joint Factor Analysis for Speaker Recognition reinterpreted as Signal Coding using Overcomplete Dictionaries", Odyssey 2010.
- [2] O. Glembek, L. Burget, N. Dehak, N. Brummer, and P. Kenny,
 "Comparison of scoring methods used in speaker recognition with Joint Factor Analysis," in *ICASSP*, 2009, pp. 4057-4060.

