



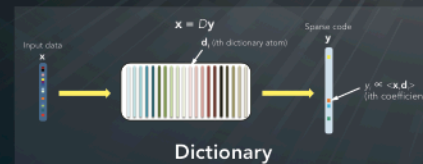
# Language Recognition via Sparse Coding

## Overview

- Language Recognition (LRE)
  - Discriminative task to identify spoken language in speech utterance
  - Acoustic signal processing techniques and learning algorithms crucial in modeling distinguishing characteristics of language
- Approach
  - Extract low-level spectral features of speech as input to sparse coding
  - Improve discriminative quality of sparse-coded speech features via maximum a posteriori (MAP) adaptation for sparse coding dictionary
- Result
  - Outperforms Lincoln i-vector pipeline developed for NIST LRE 2015 on subset comprising Arabic and Chinese clusters

## Sparse Coding Background

- Unsupervised method to learn efficient representation of data using small number of basis vectors (or atoms) from dictionary  $D$
- While solving for sparse representation of input, sparse coding also learns  $D$
- Two forms
  - $L_1$ -regularized LASSO (least absolute shrinkage and selection operator) and LARS (least angle regression)
 
$$\min_{D, y} \|x - Dy\|_2^2 + \lambda \|y\|_1 \quad \text{s.t. } \|d_i\|_2 \leq 1, \forall i$$
  - $L_0$ -regularized matching pursuit (OMP)
 
$$\min_{D, y} \|x - Dy\|_2^2 \quad \text{s.t. } \|y\|_0 \leq S$$
- Computing sparse representation of data = discovering higher-level features present in data
- Computing sparse representation of data = discovering higher-level features present in data



## Summary

- Sparse coding can also be applied to discriminative speech applications (e.g., language and speaker recognition)
- Sparse modeling can improve discriminative quality of spectral speech features
- In addition to vanilla sparse coding (VSC), we propose adaptive sparse coding (ASC), an improvement over VSC via MAP adaptation
- We experimentally validate improved performance of sparse coding over iVector-SDC on Arabic and Chinese subset from NIST LRE 2015
- Future work
  - Evaluation on full NIST LRE dataset
  - Explore different speech feature inputs such as DNN bottleneck features

## Approach

- Low-level feature extraction
  - Compute shifted delta cepstra (SDC) by windowing speech waveforms, passing through mel-scale and RASTA filterbanks, and normalizing
- Vanilla sparse coding (VSC)
  - Classical semi-supervised approach
  - Unsupervised high-level feature learning by sparse coding and dictionary learning
  - Sparse-coded features are pooled and applied to supervised training of classifiers (e.g., linear SVM)
- Adaptive sparse coding (ASC)
  - Improves VSC by adapting VSC dictionary  $D$  to utterance-specific dictionary  $D_a$  in supervised procedure
  - Using both  $D$  and  $D_a$ , compute two sparse codes
  - Arithmetic difference vector of the two sparse codes are used to train classifiers

## Experiments

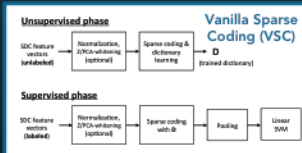
- Task: NIST Language Recognition Evaluation (LRE) 2015
  - Aims to examine average performance of language recognition system that can classify target language correctly within six predefined language clusters for a given speech sample
  - LRE 2015 covers 20 different languages in Arabic, Chinese, English, French, Slavic, and Iberian clusters
  - Our evaluation focuses on Arabic and Chinese clusters only

Cluster	Target languages
Arabic	Egyptian, Iraqi, Levantine, Maghrebi, Modern Standard
Chinese	Cantonese, Mandarin, Min, Wu

- Evaluation metric: NIST average cost performance

$$C_{avg} = \frac{1}{N} \sum_{i=1}^N C_i$$

$C_i = \frac{1}{N} \sum_{j=1}^N C_{ij}$   
 $C_{ij} = \frac{1}{N} \sum_{k=1}^N C_{ijk}$   
 $C_{ijk} = \frac{1}{N} \sum_{l=1}^N C_{ijkl}$   
 $C_{ijkl} = \frac{1}{N} \sum_{m=1}^N C_{ijklm}$   
 $C_{ijklm} = \frac{1}{N} \sum_{n=1}^N C_{ijklmn}$



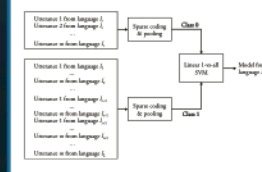
**Dictionary Adaptation Algorithm**  
Key idea: using the dictionary trained on data examples from all classes as initial point, adjust dictionary atoms to fit class-specific examples based on maximum a posteriori (MAP) criterion

**Algorithm 1** (Pseudo-code for dictionary adaptation)

1. require: arbitrary sparse model, eg.  $D_0 \in \mathbb{R}^{L \times D}$ ,  $y \in \mathbb{R}^D$ ,  $\lambda > 0$
2. initialize:  $D_1 := D_0$ ,  $A^0 := 0$ ,  $H^0 := I$
3. for  $t = 1$  to  $T$  do
  4.  $y_t := y - A^{t-1} D_0$
  5.  $D_t := \text{SparseCoding}(y_t, D_0, \lambda)$
  6.  $A^t := A^{t-1} + y_t y_t^T$
  7.  $H^t := H^{t-1} - \lambda A^t D_0 D_0^T A^t$
  8. end for
  9. return:  $D_t$



**Language Classification via SVM**  
We train 1-vs-all linear SVM classifier for each language



**Postprocessing of Classification Results**

- Historical NIST LRE systems have benefitted from fusing classification results of multiple LRE pipelines

- We use simple linear Z-score fusion based for postprocessing
- Final classification result is computed using mixing ratio  $\rho$  to combine log-likelihood ratio scores ( $llr_1$  and  $llr_2$ ) from two different classification pipelines

$$llr_{fusion} = \rho \frac{llr_1^2 - \mu_1^2}{\sigma_1^2} + (1 - \rho) \frac{llr_2^2 - \mu_2^2}{\sigma_2^2}$$

## Results: Individual Pipelines

(lower is better)

Classification pipeline	Arabic	Chinese
iVector-SDC baseline	0.2566	0.2054
Vanilla OMP-1024 sparse coding on SDC	0.2486	0.2120
Vanilla LARS-1024 sparse coding on SDC	0.2393	0.2043
Adaptive OMP-1024 sparse coding on SDC	0.2015	0.1983
Adaptive LARS-1024 sparse coding on SDC	<b>0.1874</b>	<b>0.1634</b>

- We evaluated both forms of sparse coding LARS and OMP that train 1024 dictionary atoms
- Adaptive sparse coding with LARS-1024 is our best scheme
- Significantly outperforms iVector-SDC baseline

Fusion scheme	Arabic	Chinese
iVector-SDC baseline + Vanilla LARS-1024 sparse coding on SDC	0.1988	0.1857
iVector-SDC baseline + Adaptive LARS-1024 sparse coding on SDC	0.1652	0.1226

- We calibrated mixing ratio  $\rho$  between 0.1 to 0.9
- Fusion enables us to achieve the new best cost performance for both Arabic and Chinese