

Unsupervised Speaker Diarization Using Riemannian Manifold Clustering

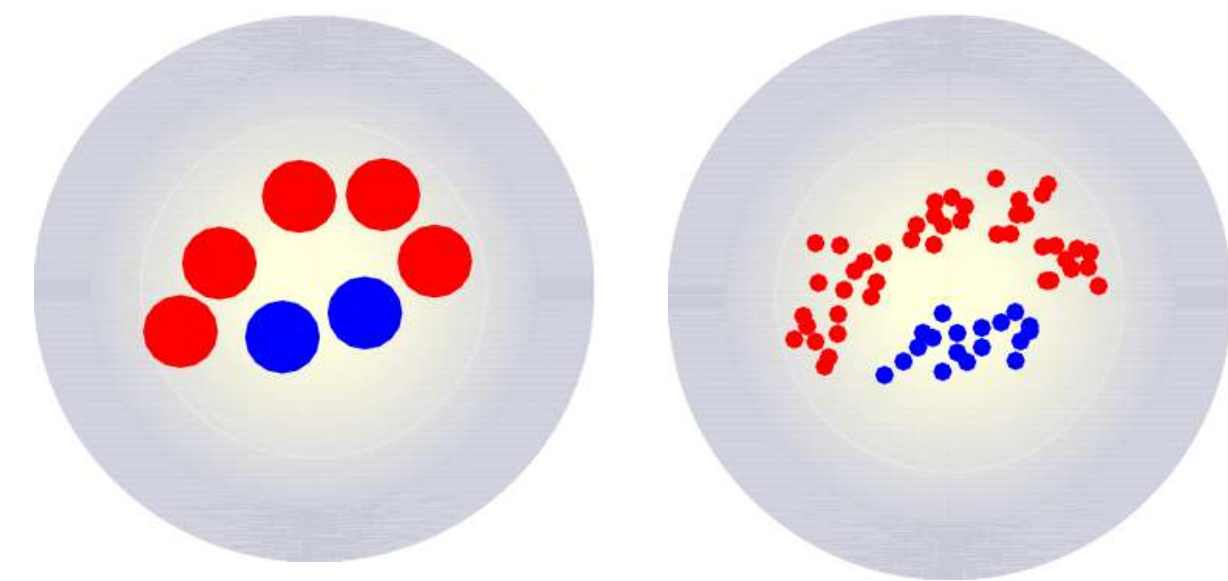
Che-Wei Huang, Bo Xiao, Panayiotis G. Georgiou, Shrikanth S. Narayanan

Motivation

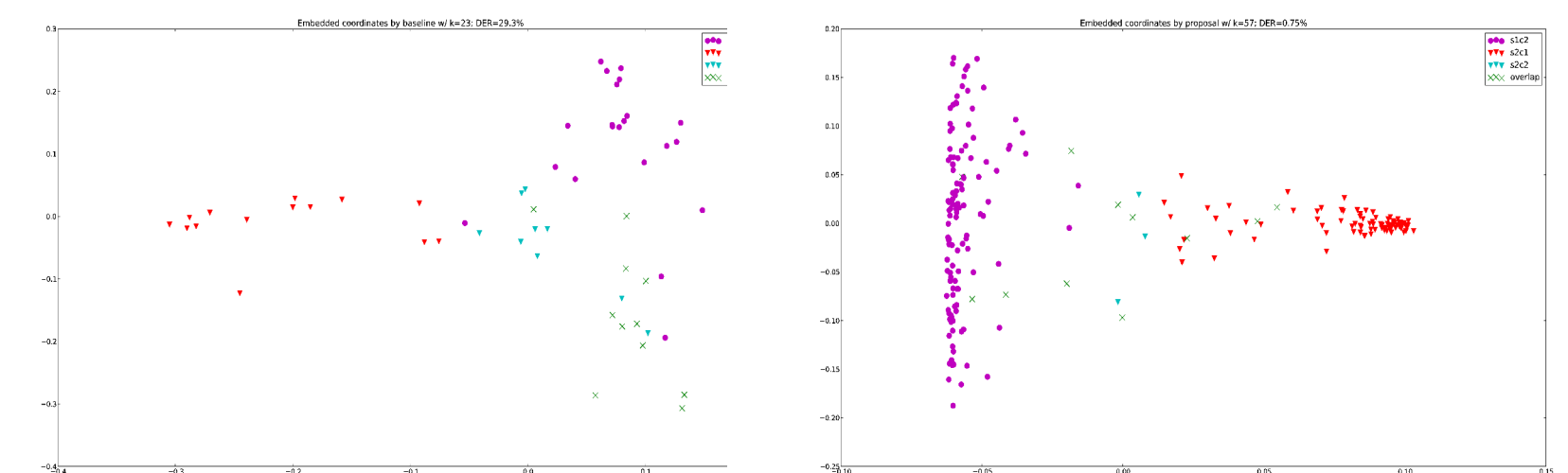
- To perfectly cluster short segments
 - Likelihood-based clustering not good on short segments
- Geometric viewpoint
 - Non-convexity of speaker clusters
 - Faithful geodesic metric
- Speaker clustering by Riemannian manifold clustering
- Suppress sparsity issue in local samples
- Stabilize performance over parameters
- Hypothesis:** speech segments from different speakers form distinct (sub-)manifolds

Proposal

- Sparsity of local samples hides clusters' structure
- Proposal:** Impose a length constraint on segments
- Advantages
 - Short segment suitable for single Gaussian
 - Increase local sample density
 - De-sensitize performance over k
- Disadvantage
 - Potential increase in computation complexity
- Schematic diagram

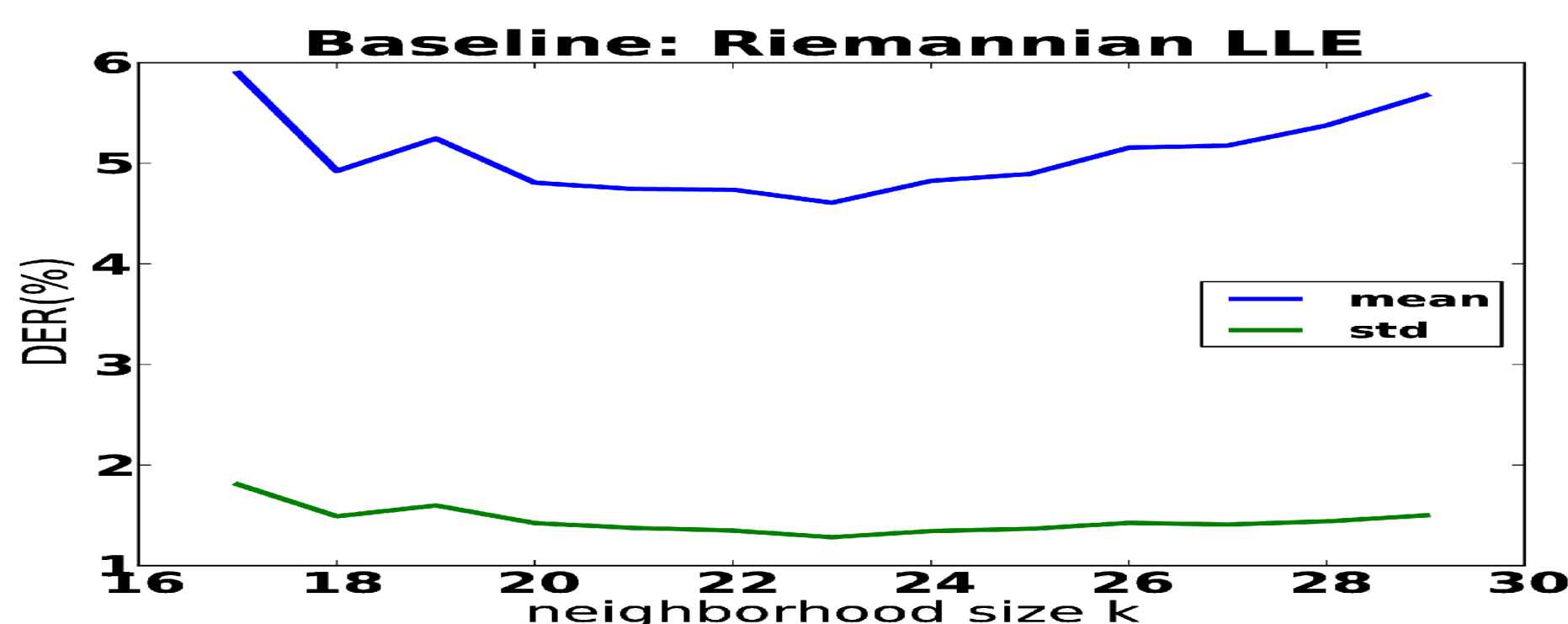


- Dense local samples reveal clusters' structure
- Embedded coordinates for a particular section



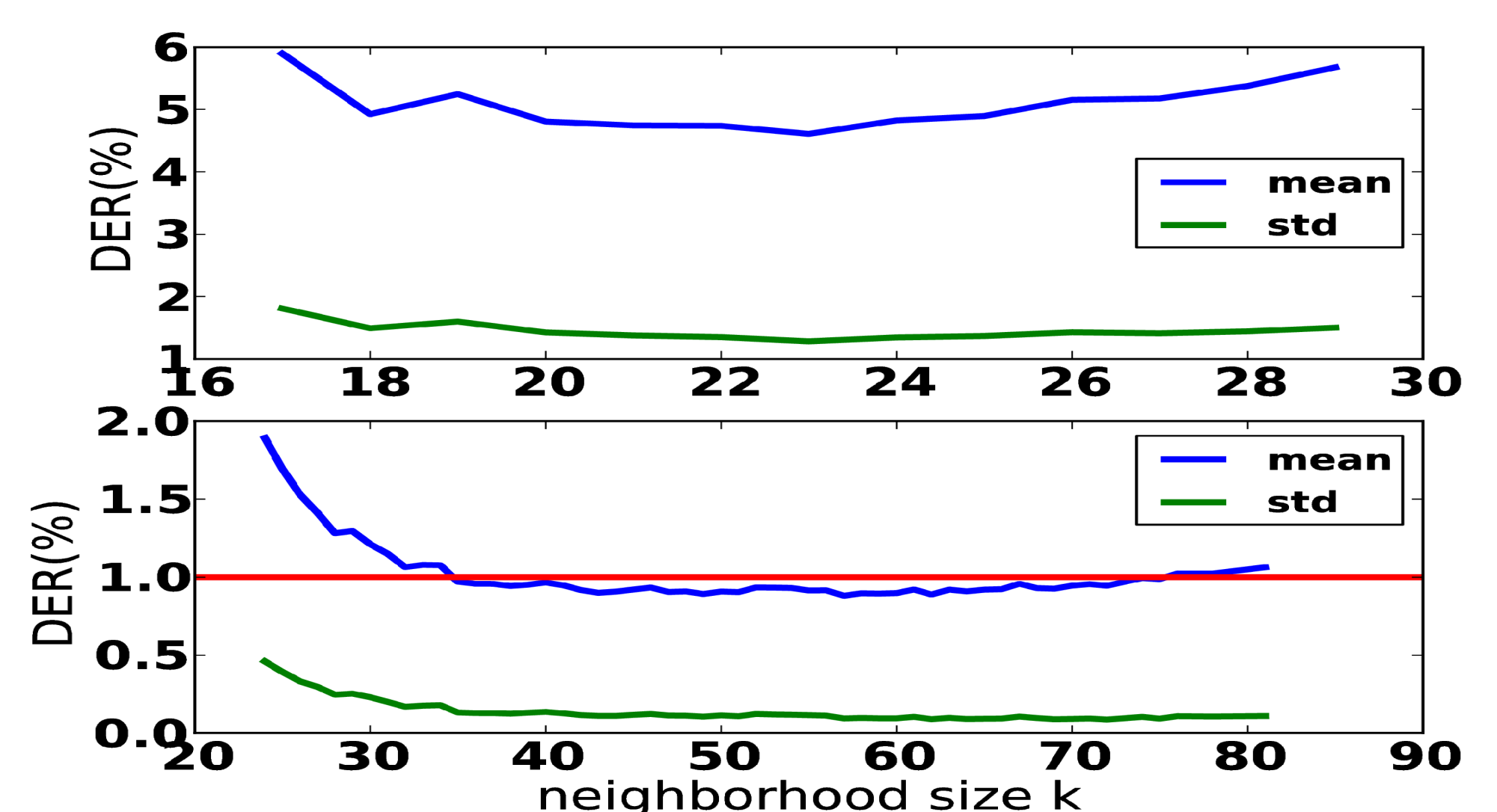
Baseline

- Riemannian Locally Linear Embedding (LLE) as the baseline
 - Each segment as a single Gaussian
 - Similarity matrix constructed by local distances to neighbors
 - Graph Laplacian for spectral clustering
- Test on 2477 5-min microphone interviews from NIST 2010 (~206hrs)



- Best at $k = 23$ w/ DER= 4.607%
- Issues
 - Choice of neighborhood size k
 - GMM vs single Gaussian modelings

Comparison



- 1s length constraint
- Best at $k = 57$ w/ DER= 0.88%
- Stable under 1% for wide range of k

Discussion & FutureWork

- Effective Riemannian manifold modeling for speech segments
- Performance less sensitive over parameter tuning at the cost of potentially higher complexity
- In the future, to consider noisy segments and the number of speakers estimation