

Sketching for Large-Scale Learning of Mixture Model

N. Keriven^{*§} A. Bourrier[†] R. Gribonval[§] P. Pérez[‡]

* Université Rennes 1, France

§ INRIA Rennes-Bretagne Atlantique, France

† Gipsa-Lab, St-Martin-d'Hères, France

‡ Technicolor, Cesson Sévigné, France

ICASSP 2016



- 1 Introduction
- 2 Proposed Algorithm
- 3 Sketching operator for Gaussian Mixture Model
- 4 Results
- 5 Conclusion

Paths to Compressive Learning

Objective

Learn parameters Θ from a **large** database $(\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^n$.

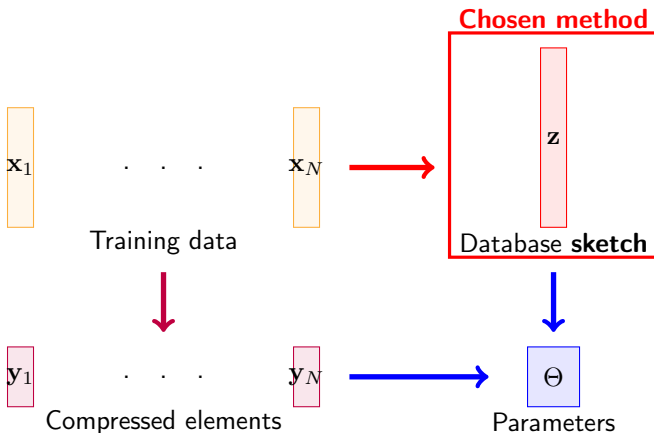
Examples:

- Learn subspace V_Θ of principal components
- Learn parameters of a classifier f_Θ
- Fit a probability distribution p_Θ
- ...

Paths to Compressive Learning

Objective

Learn parameters Θ from a **large** database $(\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^n$.



In this paper:

In this paper

Efficient method for Gaussian Mixture Model (GMM) estimation from a sketch.

Ex : Estimation of a 20-GMM from a database of 10^6 vectors in \mathbb{R}^{10}

- 5000-fold compression of the database
 - Can be performed efficiently on GPU / clusters
- Estimation process $70\times$ faster than EM
- Same precision than EM in the result

Approach : Generalized Compressive Sensing

Traditional Compressive Sensing (CS)

From $\mathbf{y} \approx \mathbf{M}\mathbf{x} \in \mathbb{R}^m$ recover vector $\mathbf{x} \in \mathbb{R}^n$

- Linear $\mathbf{M} \in \mathbb{R}^{m \times n}$ with $m < n$
- Typical assumption: \mathbf{x} sparse, etc.

Generalized Compressive Sensing

From $\mathbf{z} \approx \mathcal{A}p \in \mathbb{C}^m$ recover probability distribution $p \in L^1(\mathbb{R}^n)$

Must define:

- Linear operator $\mathcal{A} : L^1(\mathbb{R}^n) \mapsto \mathbb{C}^m$
- Generalized "sparsity" in $L^1(\mathbb{R}^n)$

Sparse probability distributions: Mixture Models

- Set of parametric probability distributions: $\mathcal{G} = \{p_{\theta}; \theta \in \mathcal{T}\}$
- " K -sparse" probability distributions :

$$p_{\Theta, \alpha} = \sum_{k=1}^K \alpha_k p_{\theta_k}$$

- Sketch $\mathbf{z} = \sum_{k=1}^K \alpha_k \mathcal{A} p_{\theta_k}$ as a combination of **atoms** in the **dictionary**:

$$\mathcal{D} = \{\mathcal{A} p_{\theta}; \theta \in \mathcal{T}\}$$

Application to Compressive Learning

Structure of the sketching operator \mathcal{A}

Collection of generalized moments $M_j : \mathbb{R}^n \mapsto \mathbb{C}$:

$$\mathcal{A}p = [\mathbb{E}_{\mathbf{x} \sim p} M_j(\mathbf{x})]_{j=1 \dots m}$$

Compressive Learning procedure

Given a database $(\mathbf{x}_1, \dots, \mathbf{x}_N) \stackrel{i.i.d.}{\sim} p$:

- Compute empirical sketch $\hat{\mathbf{z}} = [\hat{\mathbb{E}} M_j(\mathbf{x})]_{j=1 \dots m} \approx \mathcal{A}p$
- Recover $p_{\theta, \alpha}$ from $\hat{\mathbf{z}}$ using (generalized) CS techniques

Questions:

- Reconstruction algorithm ? (Section 2)
- Choice of sketching operator \mathcal{A} ? (Section 3)

Outline

- 1 Introduction
- 2 Proposed Algorithm**
- 3 Sketching operator for Gaussian Mixture Model
- 4 Results
- 5 Conclusion

Approach

Cost function

$$\min_{\Theta, \alpha} \|\hat{\mathbf{z}} - \mathcal{A}p_{\Theta, \alpha}\|_2$$

- Similar to $\min_{\mathbf{x}: \|\mathbf{x}\|_0 \leq s} \|\mathbf{y} - \mathbf{M}\mathbf{x}\|_2$ in CS.
- **Pros:** Under some hypothesis on \mathcal{G} and \mathcal{A} , yields provably good solutions with high probability (upcoming paper)
- **Cons:** Generally highly non-convex / intractable
 - Convex relaxation¹: seems difficult because of infinite / continuous dictionary
 - Greedy approaches: **approach retained here**

¹Florentina Bunea et al. **SPADES and mixture models.** *The Annals of Statistics* (2010)

Orthogonal Matching Pursuit with Replacement

- OMP: add an atom to the support by maximizing its correlation to the residual, update the residual, repeat.
- OMP with Replacement**²
 - Perform potentially **more iterations** than OMP, add a **Hard Thresholding** step.

Similar to CoSAMP or Subspace Pursuit.

- Compressive Learning OMPR** (*proposed*)
 - Enforce **non-negativity** on weights α
 - Deal with continuous dictionary using **gradient descents**
 - Add a **global optimization step** at each iteration.

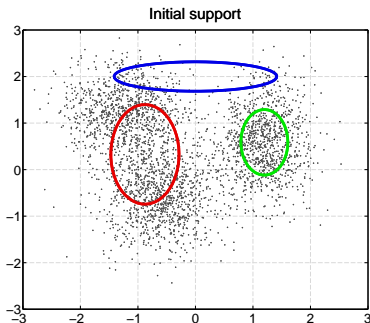
Number of iterations	Compressive Sensing	Compressive Learning
K	OMP	CLOMP
$2K$	OMPR	CLOMPR

²Prateek Jain, Ambuj Tewari, and Inderjit S. Dhillon. **Orthogonal matching pursuit with replacement.** *NIPS* (2011)

Compressive Learning OMP

Example : iteration 4 of CLOMPR, searching for a 3-GMM

- Current support

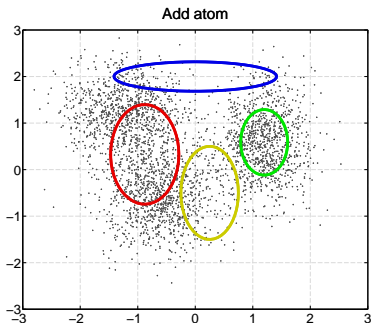


Compressive Learning OMP

Example : iteration 4 of CLOMPR, searching for a 3-GMM

- Add an atom to the support with a **gradient descent**:

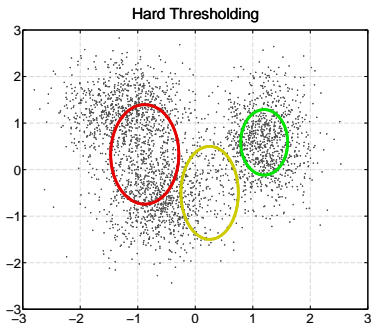
$$\arg \max_{\theta} \operatorname{Re} \left\langle \mathbf{r}, \frac{A p_{\theta}}{\|A p_{\theta}\|_2} \right\rangle$$



Compressive Learning OMP

Example : iteration 4 of CLOMPR, searching for a 3-GMM

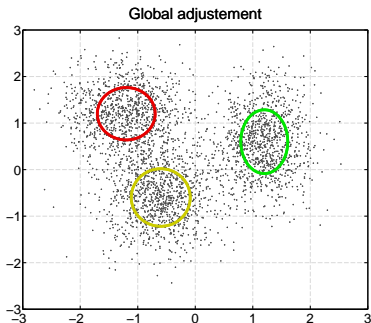
- **Hard Thresholding** to reduce the support
- Solve a **Non-negative** Least Squares to find the weights α .



Compressive Learning OMPR

Example : iteration 4 of CLOMPR, searching for a 3-GMM

- New step: **global gradient descent** initialized with the current parameters to further reduce $\|\hat{\mathbf{z}} - \mathcal{A}p_{\Theta, \alpha}\|_2$
- Update residual.



Outline

- 1 Introduction
- 2 Proposed Algorithm
- 3 Sketching operator for Gaussian Mixture Model**
- 4 Results
- 5 Conclusion

Sketching operator, Gaussian Mixture Model (GMM)

Recover $p_{\Theta, \alpha} = \sum_{k=1}^K \alpha_k p_{\theta_k}$ from $\hat{\mathbf{z}} \approx \mathcal{A}p$.

Gaussian Mixture Model

$$p_{\theta} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ with diagonal } \boldsymbol{\Sigma}$$

Random Sampling of the characteristic function³

Denote $\psi_p(\boldsymbol{\omega}) = \mathbb{E}_{\mathbf{x} \sim p}(e^{i\boldsymbol{\omega}^T \mathbf{x}})$. Given $(\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_m) \in \mathbb{R}^n$, define

$$\mathcal{A}p = [\psi_p(\boldsymbol{\omega}_j)]_{j=1, \dots, m}$$

- Analog to Random Fourier Sampling: $(\boldsymbol{\omega}_1, \dots, \boldsymbol{\omega}_m) \stackrel{i.i.d.}{\sim} \Lambda$

³Anthony Bourrier, Rémi Gribonval, and Patrick Pérez. **Compressive gaussian mixture estimation**. ICASSP (2013)

Designing the frequency distribution

The frequency distribution must "scale" with the variances of the GMM.

Approach 1 Optimize the variance of a Gaussian frequency distribution

- Requires training data with known distribution
- Classical choice⁴

⁴Dougal J Sutherland et al. **Linear-time Learning on Distributions with Approximate Kernel Embeddings.** *arXiv:1509.07553* (2015)

Designing the frequency distribution

The frequency distribution must "scale" with the variances of the GMM.

Approach 1 Optimize the variance of a Gaussian frequency distribution

Approach 2 Proposed:

- Partial preprocessing to compute the appropriate "scaling"
- Distribution that aims at maximizing $\|\nabla_{\theta} \psi_{p_{\theta}}\|_2$

The proposed distribution

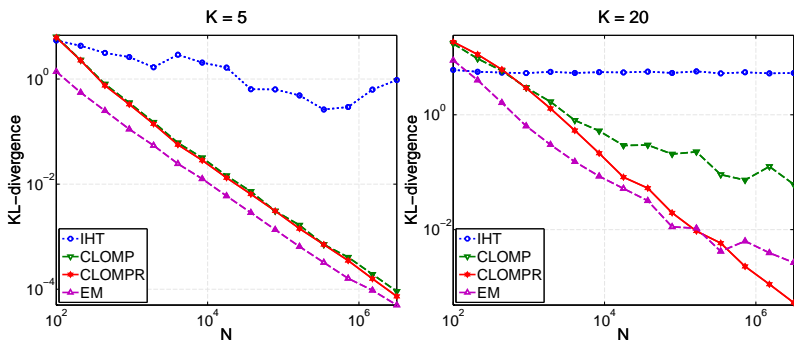
- Yields better precision in the reconstruction
- Is $20\times$ to $100\times$ faster to design

Outline

- 1 Introduction
- 2 Proposed Algorithm
- 3 Sketching operator for Gaussian Mixture Model
- 4 Results**
- 5 Conclusion

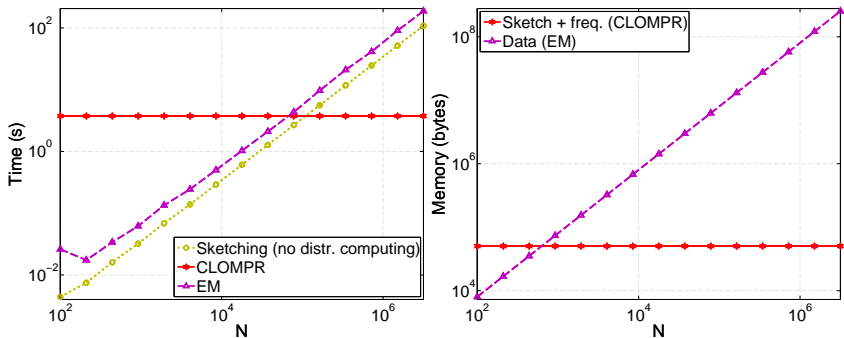
Reconstruction results

Comparison with EM (VLFeat toolbox) and previous Compressive Learning IHT⁴ (originally designed for GMM with fixed covariance). KL-div (lower is better), $n = 10$, $m = 5(2n + 1)K$.



⁴Anthony Bourrier, Rémi Gribonval, and Patrick Pérez. **Compressive gaussian mixture estimation.** *ICASSP* (2013)

Memory usage and computation time



- Sketching easily done on GPU

Application : speaker verification

- *NIST2005 database with MFCCs*
- *Classical method⁵, far from state-of-the-art but serves as a proof of concept*

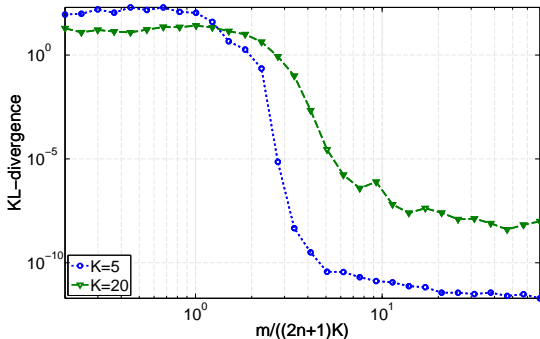
	CLOMPR			EM
	$m = 10^3$	$m = 10^4$	$m = 10^5$	
$N = 3 \cdot 10^5$	37.15	30.24	29.77	29.53
$N = 2 \cdot 10^8$	36.57	28.96	28.59	N/A

- A large database enhances the quality of the sketch

⁵Douglas A. Reynolds, Thomas F. Quatieri, and Robert B. Dunn. **Speaker Verification Using Adapted Gaussian Mixture Models.** *Digital Signal Processing* 10.1-3 (Jan. 2000)

Phase transition

Synthetic data



*Though we have preliminary theoretical guarantees (upcoming paper),
not fully explained yet.*

Outline

- 1 Introduction
- 2 Proposed Algorithm
- 3 Sketching operator for Gaussian Mixture Model
- 4 Results
- 5 Conclusion

Conclusion

Summary

Effective method to learn GMMs from a sketch, using greedy algorithms and an efficient heuristic to design the sketching operator.

Upcoming paper

- Faster algorithm for GMM with large K
- Preliminary theoretical guarantees

Future Work

- Application to other Mixture Models (α -stable distributions...)
- Generalized theoretical guarantees
- Application to other kernel methods⁶ (classification...)

⁶Dougal J Sutherland et al. **Linear-time Learning on Distributions with Approximate Kernel Embeddings.** *arXiv:1509.07553* (2015)

Questions ?

Nicolas Keriven et al. **Sketching for Large-Scale Learning of Mixture Models**. *hal-01208027v3, ICASSP (2016)*