# Sketching for Large-Scale Learning of Mixture Models

Nicolas Keriven

Université Rennes 1, Inria Rennes Bretagne-atlantique

Adv. Rémi Gribonval



#### Outline

Introduction **Practical Approach** Results Theoretical analysis 5 Conclusion and outlooks



#### **Goal** : Compute parameters $\Theta$ from a **large** database.







#### **Goal** : Compute parameters $\Theta$ from a large database.

- PCA:  $\mathbf{x} \in Span(\theta_1, ..., \theta_k)$
- Classification :  $< w_{\Theta}, \Phi(\mathbf{x}) >$
- Regression :  $\mathbf{y} = f_{\Theta}(\mathbf{x})$
- Density estimation :  $\mathbf{x} \sim p_{\Theta}$







#### **Goal** : Compute parameters $\Theta$ from a large database.

- PCA:  $\mathbf{x} \in Span(\theta_1, ..., \theta_k)$
- Classification :  $< w_{\Theta}, \Phi(\mathbf{x}) >$
- Regression :  $\mathbf{y} = f_{\Theta}(\mathbf{x})$
- Density estimation :  $\mathbf{x} \sim p_{\Theta}$

*Idea : compress the database beforehand.* 







#### **Goal** : Compute parameters $\Theta$ from a large database.





#### **Goal** : Compute parameters $\Theta$ from a large database.



• Sketch...

- Contains particular info about the database
- Maintained online
   [Cormode 2011]



- Sketch...
  Sketch...
  Maintained online [Cormode 2011]
  McContains particular info about the database
  Maintained online [Cormode 2011]
  - ...by Compressive Sensing

Recover « low-dimensional » object from few linear measurements (ex : sparse vector, low-rank matrix...)

• Sketch...

- Contains particular info about the database
- Maintained online [Cormode 2011]

• ...Learning...

Knowledge about **underlying probability distribution** 



• ...by Compressive Sensing

Recover « low-dimensional » object from few linear measurements (ex : sparse vector, low-rank matrix...)

#### Sketch = measurements of underlying probability distribution











$$\mathbf{z} = \mathcal{A}p = \left[\mathbb{E}_{\mathbf{x}\sim p}\phi_j(\mathbf{x})\right]_{j=1}^m \qquad \boldsymbol{\approx} \qquad \hat{\mathbf{z}} = \mathcal{A}\hat{p} = \left[\frac{1}{n}\sum_{i=1}^n \phi_j(\mathbf{x}_i)\right]_{j=1}^m$$



$$\mathbf{z} = \mathcal{A}p = \left[\mathbb{E}_{\mathbf{x}\sim p}\phi_j(\mathbf{x})\right]_{j=1}^m \qquad \boldsymbol{\approx} \qquad \hat{\mathbf{z}} = \mathcal{A}\hat{p} = \left[\frac{1}{n}\sum_{i=1}^n \phi_j(\mathbf{x}_i)\right]_{j=1}^m$$

Compressive Sensing : (Random) Projections





$$\mathbf{z} = \mathcal{A}p = \left[\mathbb{E}_{\mathbf{x}\sim p}\phi_j(\mathbf{x})\right]_{j=1}^m \approx \hat{\mathbf{z}} = \mathcal{A}\hat{p} = \left[\frac{1}{n}\sum_{i=1}^n \phi_j(\mathbf{x}_i)\right]_{j=1}^m$$
Compressive Sensing :  
(Random) Projections



$$\mathbf{z} = \mathcal{A}p = \left[\mathbb{E}_{\mathbf{x}\sim p}\phi_{j}(\mathbf{x})\right]_{j=1}^{m} \approx \hat{\mathbf{z}} = \mathcal{A}\hat{p} = \left[\frac{1}{n}\sum_{i=1}^{n}\phi_{j}(\mathbf{x}_{i})\right]_{j=1}^{m}$$
Compressive Sensing :  
(Random) Projections
Robustness of  
learning Alg. ?
$$\hat{\mathbf{z}} = \mathcal{A}\hat{p} = \left[\frac{1}{n}\sum_{i=1}^{n}\phi_{j}(\mathbf{x}_{i})\right]_{j=1}^{m}$$

$$\hat{\mathbf{z}} = \mathcal{A}\hat{p} = \left[\frac{1}{n}\sum_{i=1}^{n}\phi_{j}(\mathbf{x}_{i})\right]_{j=1}^{m}$$

$$\hat{\mathbf{z}} = \mathcal{A}\hat{p} = \left[\frac{1}{n}\sum_{i=1}^{n}\phi_{j}(\mathbf{x}_{i})\right]_{j=1}^{m}$$



#### Mixture Model Estimation



## Mixture Model Estimation



## Mixture Model Estimation





#### • Practical Approach (Section 2 & 3)

- Greedy algorithm inspired by Compressive Sensing
- Application to K-means, GMM with diagonal covariance

#### • Practical Approach (Section 2 & 3)

- Greedy algorithm inspired by Compressive Sensing
- Application to K-means, GMM with diagonal covariance

#### • Theoretical Analysis (Section 4)

- Information-preservation guarantee
- Infinite-dimensional Compressive Sensing



Outline

Introduction Practical Approach (Keriven, Bourrier, Gribonval, Pérèz) **Results** Theoretical analysis Conclusion and outlooks









$$\frac{M}{\mathbf{x}} \xrightarrow{\mathbf{Alg.}} \mathbf{x}_{\Gamma}$$

$$\frac{M}{\mathbf{y} = \mathbf{M}\mathbf{x}} \xrightarrow{\mathbf{Alg.}} \mathbf{x}_{\Gamma}$$

$$\mathbf{min}_{\|\mathbf{x}\|_{0} \leq s} \|\mathbf{y} - \mathbf{M}\mathbf{x}\|_{2}$$



• « Ideal » decoding scheme





- « Ideal » decoding scheme
- NP-complete



- « Ideal » decoding scheme
- NP-complete
- Two approaches:
  - Convex relaxation
  - Greedy

See [Foucart 2013]





- « Ideal » decoding scheme
- NP-complete
- Two approaches:
  - Convex relaxation
  - Greedy

See [Foucart 2013]







- « Ideal » decoding scheme
- NP-complete
- Two approaches:
  - Convex relaxation
  - Greedy



 Ideal decoding scheme (Section 4)







- « Ideal » decoding scheme
- NP-complete
- Two approaches:
  - Convex relaxation
  - Greedy



- Ideal decoding scheme (Section 4)
- Highly non-convex



See [Foucart 2013]





- « Ideal » decoding scheme
- NP-complete
- Two approaches:
  - Convex relaxation
  - Greedy



- Ideal decoding scheme (Section 4)
- Highly non-convex
- Two approaches:
  - Convex relaxation [Bunea 2010]
    - Greedy Proposed



See [Foucart 2013]

#### Proposed algorithm (Keriven, Bourrier, Gribonval, Pérèz)

#### **Orthogonal Matching Pursuit (OMP)**

[Mallat 1993, Pati 1993]

- 1. Add atom most correlated to residual
- 2. Perform Least-Squares
- 3. Repeat until desired sparsity



#### Proposed algorithm (Keriven, Bourrier, Gribonval, Pérèz)

OMP with Replacement (OMPR) [Jain 2011]

- 1. Add atom most correlated to residual
- 2. Perform Hard-Thresholding (if necessary)
- 3. Perform Least-Squares
- 4. Repeat twice desired sparsity

Similar to CoSAMP [Needell 2008] or SubSpace Pursuit [Dai 2009]



#### Proposed algorithm (Keriven, Bourrier, Gribonval, Pérèz)

#### OMP with Replacement (OMPR) [Jain 2011]

- 1. Add atom most correlated to residual
- 2. Perform Hard-Thresholding (if necessary)
- 3. Perform Least-Squares
- 4. Repeat twice desired sparsity

Similar to CoSAMP [Needell 2008] or SubSpace Pursuit [Dai 2009]

#### Compressive Learning OMPR (CLOMPR) (proposed)

- Add atom most correlated to residual with gradient descent (local min.)
- 2. Perform Hard-Thresholding
- 3. Perform Non-Negative Least-Squares
- 4. Perform gradient descent on all parameters, initialized with current ones (local min.)



We cannot just add a component

5. Repeat twice desired sparsity

## CLOMPR : illustration

(schematic illustration of one iteration)



Goal : 3-GMM. Intermediary support.



## CLOMPR : illustration

(schematic illustration of one iteration)




## CLOMPR : illustration

(schematic illustration of one iteration)







## CLOMPR : illustration

(schematic illustration of one iteration)



parameters (local min. of cost function)



## CLOMPR : illustration

(schematic illustration of one iteration)



To implement CLOMPR,  $\mathcal{A}p_{ heta}$  and  $abla_{ heta}\mathcal{A}p_{ heta}$  must have a closed-form expression w.r.t. heta



To implement CLOMPR,  $\mathcal{A}p_{ heta}$  and  $abla_{ heta}\mathcal{A}p_{ heta}$  must have a closed-form expression w.r.t. heta

#### $\boldsymbol{p}$ is spatially localized





To implement CLOMPR,  $\mathcal{A}p_{ heta}$  and  $abla_{ heta}\mathcal{A}p_{ heta}$  must have a closed-form expression w.r.t. heta

#### $p \,$ is spatially localized

#### Need incoherent sampling -> Fourier sampling









To implement CLOMPR,  $\mathcal{A}p_{ heta}$  and  $abla_{ heta}\mathcal{A}p_{ heta}$  must have a closed-form expression w.r.t. heta

#### $p\,$ is spatially localized

#### Need incoherent sampling -> Fourier sampling







$$\mathcal{A}p = \left[\psi_p(\omega_j)\right]_{j=1}^m$$

Closed-form for many models ! (including alpha-stable...)

To implement CLOMPR,  $\mathcal{A}p_{ heta}$  and  $abla_{ heta}\mathcal{A}p_{ heta}$  must have a closed-form expression w.r.t. heta

#### p is spatially localized

#### Need incoherent sampling -> Fourier sampling









#### Adjust by hand

- Not that difficult...
- The method is quite robust



#### Adjust by hand

- Not that difficult...
- The method is quite robust

#### **Cross-validation**

- Can be very long !
- Used in practice [Sutherland2015]





Adjust by hand

- Not that difficult...
- The method is quite robust

#### **Cross-validation**

- Can be very long !
- Used in practice [Sutherland2015]



- Partial pre-processing
- Heuristic based on GMMs-like distributions



### Summary



#### Given database, $\boldsymbol{m}$ , $\,\boldsymbol{K}$

- 1. Design  ${\cal A}$ 
  - Partial pre-processing to choose  $\Lambda$

• Draw 
$$(\omega_1,...,\omega_m) \overset{i.i.d.}{\sim} \Lambda$$

2. Compute 
$$\hat{\mathbf{z}} = \frac{1}{n} \left[ \sum_{i} e^{-i\omega_{j}^{T} \mathbf{x}_{i}} \right]_{j=1}^{m}$$

- Online, distributed, GPU...
- 3. Derive mixture model  $p_{\Theta, \alpha}$  with CLOMPR

Outline







#### **Classic approach**

- Goal:  $\min_{\Theta} \sum_{i=1}^{n} (\min_{1 \le l \le k} \|\mathbf{x}_i \theta_l\|_2^2)$ •
- Algorithm : Lloyd-Max [Lloyd 1982] • (Matlab's kmeans)







•



(clustered distribution = noisy mixture of Diracs)

•



Ínría-



Nicolas Keriven

### Large-scale result



 Number of measurements does not depend on N

### Large-scale result



Nicolas Keriven

# Application : spectral clustering (Keriven, Tremblay, Traonmilin, Gribonval)



n = 70 000



K-means (d=10, k=10, m=1000) Mean and var. over 50 exp.



Spectral clustering for classification [Uw 2001], augmented MNIST database [Loosli 2007].



# Application : spectral clustering (Keriven, Tremblay, Traonmilin, Gribonval)



K-means (d=10, k=10, m=1000) Mean and var. over 50 exp.



Spectral clustering for classification [Uw 2001], augmented MNIST database [Loosli 2007].



# Application : spectral clustering (Keriven, Tremblay, Traonmilin, Gribonval)



K-means (d=10, k=10, m=1000) Mean and var. over 50 exp.



Spectral clustering for classification [Uw 2001], augmented MNIST database [Loosli 2007].

 CLOMPR performs better and is more stable with a large database

# Application : speaker recognition

Variant of CLOMPR, faster at large k		R, (Hierai	(Hierarchical) CLOMPR		
		$m = 10^3$	$m = 10^4$	$m = 10^5$	
	$n = 3.10^5$	37.15	30.24	29.77	29.53
	$n = 2.10^{8}$	36.57	28.96	28.59	N/A

Classical method for speaker recognition [Reynolds 2000] (for proof of concept) NIST 2005 database, MFCCs.

• Also performs better on a large database.

### Outline

Introduction **Practical Approach** Results

**Theoretical analysis** (Gribonval, Blanchard, Keriven, Traonmilin) **Conclusion and outlooks** 



#### Information-preservation guarantees

Guarantee for CLOMPR ? Difficult ! (non-convex, random...)



### Information-preservation guarantees

#### Guarantee for CLOMPR ? Difficult ! (non-convex, random...)





### Information-preservation guarantees



- Robustness to using  $\hat{\mathbf{z}} = \mathcal{A}\hat{p}~$  instead of  $\mathbf{z} = \mathcal{A}p$  ?
- Robustness to p not being **exactly** a mixture model ?
- Guarantees in terms of usual learning cost functions ?
  - K-means : sum of distances to closest centroid
  - GMMs : negative log-likelihood

### K-means: result

**Goal** minimize 
$$R(\Theta) = \mathbb{E}_{\mathbf{x} \sim p^*} \left[ \min_l \|\mathbf{x} - \theta_l\|_2^2 \right]$$
 (expected risk)



### K-means: result

**Goal** minimize 
$$R(\Theta) = \mathbb{E}_{\mathbf{x} \sim p^*} \left[ \min_l \|\mathbf{x} - \theta_l\|_2^2 \right]$$
 (expected risk)



•  $\mathcal{E}$  - separation

- M bounded domain
- **Reweighted** Fourier features (needed for theory, no effect in practice)

### K-means: result

**Goal** minimize 
$$R(\Theta) = \mathbb{E}_{\mathbf{x} \sim p^*} \left[ \min_l \|\mathbf{x} - \theta_l\|_2^2 \right]$$
 (expected risk)



- on • M-bo
  - M bounded domain
  - **Reweighted** Fourier features (needed for theory, no effect in practice)

If 
$$m \ge O\left(k^2 d^3 \text{polylog}(k, d) \log(M/\varepsilon)\right)$$
  
w.h.p.  $R(\hat{\Theta}) \lesssim R(\Theta^*) + O\left(\sqrt{d^2 k/n}\right)$   
Minimize cost func.  
(with hyp.) Minimize expected risk  
(with hyp.) Nicolas Keriven 17/28

## GMMs with known covariance : result

**Goal** minimize 
$$R(\Theta, \alpha) = \mathbb{E}_{\mathbf{x} \sim p^*} \left[ -\log p_{\Theta, \alpha}(\mathbf{x}) \right]_{\text{(expected risk)}}^{p_{\Theta, \alpha} = \sum_l \alpha_l \mathcal{N}(\theta_l, \Sigma)}$$



### GMMs with known covariance : result

**Goal** minimize 
$$R(\Theta, \alpha) = \mathbb{E}_{\mathbf{x} \sim p^*} \left[ -\log p_{\Theta, \alpha}(\mathbf{x}) \right]_{\text{(expected risk)}}^{p_{\Theta, \alpha} = \sum_l \alpha_l \mathcal{N}(\theta_l, \Sigma)}$$

#### Large enough separation





- M bounded domain
- Fourier features

### GMMs with known covariance : result

**Goal** minimize 
$$R(\Theta, \alpha) = \mathbb{E}_{\mathbf{x} \sim p^*} \left[ -\log p_{\Theta, \alpha}(\mathbf{x}) \right]_{\text{(expected risk)}}^{p_{\Theta, \alpha} = \sum_l \alpha_l \mathcal{N}(\theta_l, \Sigma)}$$

#### Large enough separation





- M bounded domain
- Fourier features

Ifm large enoughw.h.p.
$$R(\hat{\Theta}, \hat{\alpha}) - R(\Theta^*, \alpha^*) \lesssim \inf_{\Theta, \alpha} \|p^* - p_{\Theta, \alpha}\|_{L^1} + \mathcal{O}(1/\sqrt{n})$$
 $\checkmark$ L1 distance from p\* to the set of (separated) GMMs

### GMM trade-off

More high



### GMM with **unknown diagonal** covariance


## GMM with **unknown diagonal** covariance





### Number of measurements

In theory, at least

 $m \geq \mathcal{O}(k^2 d^2)$ 



### Number of measurements

Relative **K-means** In theory, at least SSE 0.02 0.04 0.06 0.08 0.1 0 0.02 0.04 0.06  $m \ge \mathcal{O}(k^2 d^2)$ 30 30 25 25 20 20 **Empirically** ? σ × 15 15  $m \approx \mathcal{O}(kd)$ 10 10 5 5  $10^{0}$ 10<sup>-1</sup> 10<sup>1</sup> 10<sup>-1</sup>  $10^{0}$  $10^{1}$ m/(kd)m/(kd)Relative Relative GMMs, diagonal cov. GMMs, known cov. loglike loglike >2 >2 1.5 >2 1.5 1.5 1 >2 1 1.5 1 30 30 30 30 25 25 25 25 20 20 20 20 × σ × σ 15 15 15 15 10 10 10 10 5 5 5 5 10<sup>0</sup> 10<sup>0</sup> 10<sup>0</sup> 10<sup>0</sup> 10<sup>1</sup> 10<sup>-1</sup> 10<sup>1</sup>  $10^{-1}$ 10<sup>1</sup> 10<sup>-1</sup>  $10^{1}$  $10^{-1}$ m/(kd)m/(kd) m/(kd) m/(kd)

# Sketch of proof : principle

Goal : Existence of instance Optimal Decoder



$$\|p^* - \Delta(\hat{\mathbf{z}})\| \lesssim d(p^*, \mathfrak{S}) + \underbrace{\|\mathcal{A}(p^* - \hat{p})\|}_{\mathcal{O}(1/\sqrt{n})}$$

# Sketch of proof : principle





# Sketch of proof : principle



1: Proving non-uniform LRIP





1: Proving non-uniform LRIP



Kernel mean embedding [Smola 2007] Random (Fourier) Features [Rahimi 2007]

$$\|\mathcal{A}(q-q')\|_{2}^{2} \approx \|q-q'\|_{\kappa}^{2}$$

Hoeffding, Bernstein, chaining...













### Results

#### **Sufficient** Conditions

• S has finite covering numbers

 $\eta = \mathcal{O}(1/\sqrt{m})$ 

*Ex* : *GMMs* with unknown covariance

Not great !

### Results

#### **Sufficient** Conditions

- S has finite covering numbers
- S mixtures of sufficiently separated distributions
- $\kappa(p_{\theta}, p_{\theta'}) = f(\|\theta \theta'\|)$ with smooth f
- « Smooth » Random Features
- Smooth risk  ${\cal R}$

$$\eta = \mathcal{O}(1/\sqrt{m})$$

Not great !

#### *Ex* : *GMMs* with unknown covariance

$$\eta = \mathcal{O}(C^{-m})$$

+ guarantees w.r.t. risk

Ex : Mixture of Diracs (K-means) with  $m \geq \mathcal{O}\left(k^2 d^2 \text{polylog}(k, d) \log(1/\eta)\right)$ 



### Results

#### **Sufficient** Conditions

- S has finite covering numbers
- S mixtures of sufficiently separated distributions
- $\kappa(p_{\theta}, p_{\theta'}) = f(\|\theta \theta'\|)$ with smooth f
- « Smooth » Random Features
- Smooth risk  ${\cal R}$
- « Smoother » Random Features

 $\eta = \mathcal{O}(1/\sqrt{m})$ 

Not great !

#### *Ex : GMMs with unknown covariance*

$$\eta = \mathcal{O}(C^{-m})$$

#### + guarantees w.r.t. risk

Ex : Mixture of Diracs (K-means) with  $m \geq \mathcal{O}\left(k^2 d^2 \text{polylog}(k, d) \log(1/\eta)\right)$ 

 $\eta = 0$ 

#### *Ex :*

- Mixtures of Diracs (K-means) with  $m \geq \mathcal{O}\left(k^2 \mathbf{d}^3 \operatorname{polylog}(k, d)\right)$
- GMMs with known covariance

# Outline

Introduction **Practical Approach** Results Theoretical analysis **Conclusion and outlooks** 



Greedy algorithm for large-scale mixture learning from random moments



- Greedy algorithm for large-scale mixture learning from random moments
- Efficient heuristic to design the sketching operator as Fourier sampling



- Greedy algorithm for large-scale mixture learning from random moments
- Efficient heuristic to design the sketching operator as Fourier sampling
- Application to **mixtures of Diracs, GMMs**



- Greedy algorithm for large-scale mixture learning from random moments
- Efficient heuristic to design the sketching operator as Fourier sampling
- Application to **mixtures of Diracs, GMMs**
- Evaluation on **synthetic** and **real** data



- Greedy algorithm for large-scale mixture learning from random moments
- Efficient heuristic to design the sketching operator as Fourier sampling
- Application to mixtures of Diracs, GMMs
- Evaluation on **synthetic** and **real** data
- Information preservation guarantees using infinite-dimensional Compressive Sensing



# The SketchMLbox

### SketchMLbox (sketchml.gforge.inria.fr)

- Mixture of Diracs (« K-means »)
- GMMs with known covariance
- GMMs with unknown diagonal covariance
- Soon:
  - Alpha-stable
  - Gaussian Locally Linear Mapping [Deleforge 2014]
- Optimized for user-defined  $(Ap_{\theta}, \nabla_{\theta}Ap_{\theta})$









Algorithmic guarantees ? Non-convex cost function, randomized algorithm...



Algorithmic guarantees ? Non-convex cost function, randomized algorithm...

• Locally convex ?



Algorithmic guarantees ? Non-convex cost function, randomized algorithm...

- Locally convex ?
- Basin of attraction ? [Jacques 2016, Candes 2016...]





#### Algorithmic guarantees ? Non-convex cost function, randomized algorithm...

- Locally convex ?
- Basin of attraction ? [Jacques 2016, Candes 2016...]



Reached by CLOMPR with reasonable hypotheses ?



#### Algorithmic guarantees ? Non-convex cost function, randomized algorithm...

- Locally convex ?
- Basin of attraction ? [Jacques 2016, Candes 2016...]



- Reached by CLOMPR with reasonable hypotheses ?
- Stopping condition ?





#### Algorithmic guarantees ? Non-convex cost function, randomized algorithm...

- Locally convex ?
- Basin of attraction ? [Jacques 2016, Candes 2016...]
- $\begin{array}{c} f(z) & -f(z) \\ 1 \\ 2 \\ 1 \\ 0 \\ -2 \\ -2 \\ -1 \\ 0 \\ -2 \\ -2 \\ -1 \\ 0 \\ -2 \\ -2 \\ -1 \\ 0 \\ -2 \\ -2 \\ -1 \\ 0 \\ -2 \\ -2 \\ -1 \\ 0 \\ -2 \\ -2 \\ -1 \\ 0 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\ -1 \\ 0 \\ -2 \\$
- Reached by CLOMPR with reasonable hypotheses ?
- Stopping condition ?

#### Recent result : locally block convex







1. Bridge observed gap between theory and practice ?





- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from  $\mathcal{E}$  coverings





- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from  $\mathcal{E}$  coverings
  - Improve concentration inequalities ?





- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from  $\mathcal{E}$  coverings
  - Improve concentration inequalities ?
- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?





- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from  $\mathcal{E}$  coverings
  - Improve concentration inequalities ?
  - 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
  - 3. Extend framework to other tasks ?





- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from  $\mathcal{E}$  coverings
  - Improve concentration inequalities ?
- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
- 3. Extend framework to other tasks ?
  - Recent paper submitted to AISTATS : **PCA**





- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from *C* coverings
  - Improve concentration inequalities ?



- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
- 3. Extend framework to other tasks ?
  - Recent paper submitted to AISTATS : **PCA**
  - Other existing use of Fourier sketches ? : e.g. classification [Sutherland 2015]



- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from *C* coverings
  - Improve concentration inequalities ?



- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
- 3. Extend framework to other tasks ?
  - Recent paper submitted to AISTATS : **PCA**
  - Other existing use of Fourier sketches ? : e.g.  $cla: K(\mathbf{M}, \mathbf{M}) \approx z(\mathbf{M})^T z(\mathbf{M})$
  - Other kernel methods (algorithmic ? Theoretical ?)

- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from *C* coverings
  - Improve concentration inequalities ?



- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
- 3. Extend framework to other tasks ?
  - Recent paper submitted to AISTATS : **PCA**
  - Other existing use of Fourier sketches ? : e.g.  $cla: K(\mathbf{M}, \mathbf{M}) \approx z(\mathbf{M})^T z(\mathbf{M})$
  - Other kernel methods (algorithmic ? Theoretical ?)
- 4. Extension to multi-layer sketches ? (Neural networks...)


# Outlooks : extension of the methods

- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from *C* coverings
  - Improve concentration inequalities ?



- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
- 3. Extend framework to other tasks ?
  - Recent paper submitted to AISTATS : **PCA**
  - Other existing use of Fourier sketches ? : e.g.  $cla: K(\mathbf{M}, \mathbf{M}) \approx z(\mathbf{M})^T z(\mathbf{M})$
  - Other kernel methods (algorithmic ? Theoretical ?)
- 4. Extension to multi-layer sketches ? (Neural networks...)
  - May be adapted to e.g. GMMs with unknown covariance



# Outlooks : extension of the methods

- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from *C* coverings
  - Improve concentration inequalities ?



- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
- 3. Extend framework to other tasks ?
  - Recent paper submitted to AISTATS : **PCA**
  - Other existing use of Fourier sketches ? : e.g.  $cla: K(\mathbf{M}, \mathbf{M}) \approx z(\mathbf{M})^T z(\mathbf{M})$
  - Other kernel methods (algorithmic ? Theoretical ?)
- 4. Extension to multi-layer sketches ? (Neural networks...)
  - May be adapted to e.g. GMMs with unknown covariance
  - Equivalence between LRIP and instance optimality still valid for non-linear operators !

## Outlooks : extension of the methods

- 1. Bridge observed gap between theory and practice ?
  - Does *not* come from @verings
  - Improve concentration inequalities ?



 $\mathsf{K}(\mathbf{\overline{4}},\mathbf{\overline{4}}) \approx \mathsf{Z}(\mathbf{\overline{4}})^{\mathsf{T}}\mathsf{Z}(\mathbf{\overline{4}})$ 

- 2. Combine with « regular » dimensionality reduction for both tall and fat databases ?
- 3. Extend framework to other tasks ?
  - Recent paper submitted to AISTATS : PCA
  - Other existing use of Fourier sketches ? : e.g. classification [Sutherland 2015]
  - Other kernel methods (algorithmic ? Theoretical ?)
- 4. Extension to multi-layer sketches ? (Neural networks...)
  - May be adapted to e.g. GMMs with unknown covariance
  - Equivalence between LRIP and instance optimality still valid for non-linear operators
    !
  - CLOMPR and current sufficient conditions no longer valid...



# Thank you !

- K., Bourrier, Gribonval, Perez. Sketching for Large-Scale Learning of Mixture Models *ICASSP* 2016
- K., Bourrier, Gribonval, Perez. Sketching for Large-Scale Learning of Mixture Models (extended version) *submitted to Information and Inference, arXiv:1606.0238*
- K., Tremblay, Gribonval, Traonmilin. Compressive K-means ICASSP 2017
- K., Tremblay, Gribonval. SketchMLbox (sketchml.gforge.inria.fr)
- Gribonval, Blanchard, K., Traonmilin. Random moments for Sketched Statistical Learning submitted to AISTATS 2017, extended version soon





Nicolas Keriven

# Appendix : CLOMPR

Algorithm 2: Compressive mixture learning à la OMP: CLOMP (T = K) and CLOMPR (T = 2K)**Data**: Empirical sketch  $\hat{\mathbf{z}}$ , sketching operator  $\mathcal{A}$ , sparsity K, number of iterations T > K**Result**: Support  $\Theta$ , weights  $\alpha$  $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{z}}; \Theta \leftarrow \emptyset;$ for  $t \leftarrow 1$  to T do Step 1: Find a normalized atom highly correlated with the residual with a gradient descent  $\theta \leftarrow \text{maximize}_{\theta} \left( \operatorname{Re} \left\langle \frac{AP_{\theta}}{\|AP_{\theta}\|_{2}}, \hat{\mathbf{r}} \right\rangle_{2}, \text{init} = \text{rand} \right);$ end Step 2: Expand support  $\Theta \leftarrow \Theta \cup \{\theta\};$ end Step 3: Enforce sparsity by Hard Thresholding if needed if  $|\Theta| > K$  then  $\boldsymbol{\beta} \leftarrow \arg\min_{\boldsymbol{\beta} \ge 0} \left\| \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \beta_k \frac{\mathcal{A} P_{\boldsymbol{\theta}_k}}{\left\| \mathcal{A} P_{\boldsymbol{\theta}_k} \right\|_2} \right\|_2 \text{ Select } K \text{ largest entries } \beta_{i_1}, \dots, \beta_{i_K};$ Reduce the support  $\Theta \leftarrow \{\theta_{i_1}, ..., \theta_{i_K}\};$ end end Step 4: Project to find weights  $\alpha \leftarrow \arg\min_{\alpha \ge 0} \left\| \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \alpha_k \mathcal{A} P_{\boldsymbol{\theta}_k} \right\|_{2};$ end **Step 5**: Perform a gradient descent *initialized with current parameters*  $\Theta, \alpha \leftarrow \texttt{minimize}_{\Theta, \alpha} \left( \left\| \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \alpha_k \mathcal{A} P_{\boldsymbol{\theta}_k} \right\|_{\alpha}, \texttt{init} = (\Theta, \alpha), \texttt{constraint} = \{ \alpha \ge 0 \} \right);$ end Update residual:  $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \alpha_k \mathcal{A} P_{\boldsymbol{\theta}_k}$ end Normalize  $\alpha$  such that  $\sum_{k=1}^{K} \alpha_k = 1$