# Sketched Learning from Random Features Moments

**Nicolas Keriven**
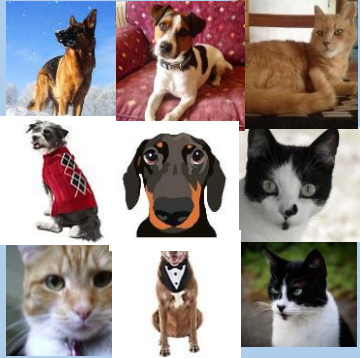
Ecole Normale Supérieure (Paris)

CFM-ENS chair in Data Science

*(thesis with Rémi Gribonval at Inria Rennes)*

Imaging in Paris, Apr. 5th 2018
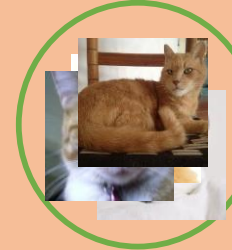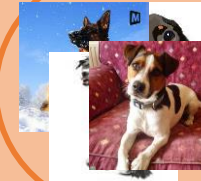
**Database**

**Learning**

**Task**

- Clustering

- Classification

= cat

- etc…

**Large** database

**Task**

Learning

**Slow, costly**

- Clustering

- Classification

= cat

- etc…

**Large database**

**Distributed database**

Learning

*Slow, costly*

**Task**

- Clustering

- Classification

= cat

- etc...

**Large database**

**Learning**

**Slow, costly**

**Task**

- Clustering

- Classification

= cat

- etc…

**Distributed database**

**Data Stream**

…        …

# Context: machine learning



**Large** database

**Distributed** database

Data **Stream**

Learning

**Slow, costly**

**Task**

- Clustering

- Classification

= cat

- etc...

**Idea!**

**Small** intermediate representation

**Large** database

**Learning**

**Slow, costly**

**Task**

- Clustering

- Classification

= cat

- etc...

**Distributed** database

**Data Stream**

...     ...

**1: Compression**

**Idea!**

**Small** intermediate representation

**Large database**

**Task**

Learning

*Slow, costly*

- Clustering

- Classification

= cat

- etc...

**Distributed database**

**Data *Stream***

**1: Compression**

**2: Learning**

*Idea!*

**Small** intermediate representation

**Large database**

**Distributed database**

**Data Stream**

...   ...

**Learning**

*Slow, costly*

**Task**

- Clustering

- Classification

= cat

- etc...

**2: Learning**

*Idea!*

**1: Compression**

**Small intermediate representation**

**Desired properties**
- **Fast** to compute (distributed, streaming, **GPU**...)
- Preserve desired **information**
- Preserve **data privacy**

**Database**

*Feature extraction*

$$n$$

$$d \quad x_1 | x_2 \quad \cdot \ \cdot \ \cdot \quad x_n$$

Data = Collection of vectors

**Database**



*Feature extraction*

$\longrightarrow$

$n$

$d$

$x_1$ $x_2$ . . . $x_n$

*Compression ?*

Data = Collection of vectors

# Three compression schemes

**Database**



Feature extraction →

$n$

$d$ | $x_1$ | $x_2$ | . . . | $x_n$

Data = Collection of vectors

Compression ?

$n$

$d'$ | $x'_1$ | $x'_2$ | . . . | $x'_n$

**Dimensionality reduction**

See eg *[Calderbank 2009, Boutsidis 2010]*

- Random Projection
- Feature selection

# Three compression schemes

**Database**



Feature extraction →

$n$

$d$ | $x_1$ $x_2$ . . . $x_n$

Compression ?

Data = Collection of vectors

---

$n$

$d'$ | $x'_1$ $x'_2$ . . . $x'_n$

**Dimensionality reduction**
See eg *[Calderbank 2009, Boutsidis 2010]*

- Random Projection
- Feature selection

---

**Subsampling coresets**
See eg
*[Feldman 2010]*

$n'$

$d$ | $x_1$ . . . $x_{n'}$

- Uniform sampling (naive)
- Adaptive sampling…

# Three compression schemes

**Database**



*Feature extraction* $\rightarrow$

$d$ — $x_1$ $x_2$ . . . $x_n$ — $n$

*Compression ?*

Data = Collection of vectors

---

$n$

$d'$ — $x'_1$ $x'_2$ . . . $x'_n$

**Dimensionality reduction**

See eg *[Calderbank 2009, Boutsidis 2010]*

- Random Projection
- Feature selection

---

**Subsampling coresets**

See eg *[Feldman 2010]*

$n'$

$d$ — $x_1$ . . . $x_{n'}$

- Uniform sampling (naive)
- Adaptive sampling…

---

**Linear sketch**

See *[Thaper 2002] [Cormode 2011]*

*Distributed, streaming*

$m$ — **z**

- Hash tables, histograms
- **Sketching for learning ?**

**What is a sketch ?**

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

**What is a sketch ?**

Any *linear* sketch = **empirical moments**

$$\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X) = \frac{1}{n}\sum_i \Phi(x_i)$$

**What is contained in a sketch ?**

**What is a sketch ?**

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

**What is contained in a sketch ?**

- $\Phi(x) = x$ : mean

**What is a sketch ?**

Any *linear* sketch = **empirical moments**

$$\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X) = \frac{1}{n}\sum_i \Phi(x_i)$$

**What is contained in a sketch ?**

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

### What is a sketch ?

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

### What is contained in a sketch ?

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

**What is a sketch ?**

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

**What is contained in a sketch ?**

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

- Proposed: **kernel random features**
  *[Rahimi 2007]*
  (random proj. + non-linearity)

## What is a sketch ?

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

## What is contained in a sketch ?

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

- Proposed: **kernel random features**
  *[Rahimi 2007]*
  (random proj. + non-linearity)

## Questions:

- **What information is preserved by the sketching ?**

- How to retrieve this information ?

- **What is a sufficient number of features ?**

## What is a sketch ?

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

## What is contained in a sketch ?

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\mathrm{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

- Proposed: **kernel random features**
  *[Rahimi 2007]*
  (random proj. + non-linearity)

## Questions:

- **What information is preserved by the sketching ?**

- How to retrieve this information ?

- **What is a sufficient number of features ?**

**Intuition: sketching as a linear embedding**

## What is a sketch ?

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

## What is contained in a sketch ?

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

- Proposed: **kernel random features**
  *[Rahimi 2007]*
  (random proj. + non-linearity)

## Questions:

- **What information is preserved by the sketching ?**

- How to retrieve this information ?

- **What is a sufficient number of features ?**

**Intuition: sketching as a linear embedding**

- Assumption: $\quad x_1, ..., x_n \overset{i.i.d.}{\sim} \pi^\star$

## What is a sketch ?

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

## What is contained in a sketch ?

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

- Proposed: **kernel random features**
  *[Rahimi 2007]*
  (random proj. + non-linearity)

## Questions:

- **What information is preserved by the sketching ?**

- How to retrieve this information ?

- **What is a sufficient number of features ?**

## Intuition: sketching as a **linear embedding**

- Assumption: $x_1, ..., x_n \overset{i.i.d.}{\sim} \pi^\star$

- Linear operator: $\mathcal{A}\pi = \mathbb{E}_{X \sim \pi}\Phi(X)$

## What is a sketch ?

Any *linear* sketch = **empirical moments**

$$\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X) = \frac{1}{n}\sum_i \Phi(x_i)$$

## What is contained in a sketch ?

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

- Proposed: **kernel random features**
  *[Rahimi 2007]*
  (random proj. + non-linearity)

## Questions:

- **What information is preserved by the sketching ?**

- How to retrieve this information ?

- **What is a sufficient number of features ?**

## Intuition: sketching as a **linear embedding**

- Assumption: $x_1, ..., x_n \overset{i.i.d.}{\sim} \pi^\star$

- Linear operator: $\mathcal{A}\pi = \mathbb{E}_{X \sim \pi}\Phi(X)$

- « Noisy » linear measurement:

$$\hat{\mathbf{z}} = \mathcal{A}\pi^\star + \hat{\mathbf{e}}$$

*Noise* $\hat{\mathbf{e}} = \hat{\mathbb{E}}\Phi(X) - \mathbb{E}_{\pi^\star}\Phi(X)$ *small*

## What is a sketch ?

Any *linear* sketch = **empirical moments**

$$\boxed{\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)} = \frac{1}{n}\sum_i \Phi(x_i)$$

## What is contained in a sketch ?

- $\Phi(x) = x$ : mean

- $\Phi(x) = x^k$ : $k^{\text{th}}$ moment

- $\Phi(x) = [1_{x \in B_i}]_{i=1}^m$ : histogram

- Proposed: **kernel random features**
  *[Rahimi 2007]*
  (random proj. + non-linearity)

## Questions:

- **What information is preserved by the sketching ?**

- How to retrieve this information ?

- **What is a sufficient number of features ?**

## Intuition: sketching as a linear embedding

- Assumption: $x_1, ..., x_n \overset{i.i.d.}{\sim} \pi^\star$

- Linear operator: $\mathcal{A}\pi = \mathbb{E}_{X \sim \pi}\Phi(X)$

- « Noisy » linear measurement:

$$\boxed{\hat{\mathbf{z}} = \mathcal{A}\pi^\star + \hat{\mathbf{e}}}$$

*Noise* $\hat{\mathbf{e}} = \hat{\mathbb{E}}\Phi(X) - \mathbb{E}_{\pi^\star}\Phi(X)$ *small*

**Dimensionality-reducing, random, linear embedding: Compressive Sensing?**

**Compressive Sensing:**

**Classical compressive sensing**

$$\mathbf{x}$$

**Sketched learning in this talk**

$$\pi^\star$$

**Compressive Sensing:**

- Dimensionality reduction, random operator

**Classical compressive sensing**

*Random matrix*

$$\mathbf{y} = \mathbf{M}\mathbf{x} + \mathbf{e}$$

$\mathbf{x}$

**Sketched learning in this talk**

*Random features averaged*

$$\mathbf{z} = \mathcal{A}\pi^\star + \mathbf{e}$$

$\pi^\star$

# Compressive Sensing: sparsity ?

## Compressive Sensing:

- Dimensionality reduction, random operator

- (Ill-posed) **inverse problem**: *density estimation*

## Classical compressive sensing



*Random matrix*

$$\mathbf{y} = \mathbf{M}\mathbf{x} + \mathbf{e}$$

$\mathbf{x}$

## Sketched learning in this talk



*Random features averaged*

$$\mathbf{z} = \mathcal{A}\pi^{\star} + \mathbf{e}$$

$\pi^{\star}$

**Compressive Sensing:**

- Dimensionality reduction, random operator

- (Ill-posed) **inverse problem**: *density estimation*

- **Sparsity**: « simple » densities (mixture model)

**Classical compressive sensing**



$$\mathbf{y} = \mathbf{Mx} + \mathbf{e}$$

$\mathbf{x}$     $\mathbf{x}_k$

**Sketched learning in this talk**



$$\mathbf{z} = \mathcal{A}\pi^{\star} + \mathbf{e}$$

$\pi^{\star}$     $\pi_{\text{mix.}}$

**Mixture of Diracs = k-means**

**Mixture of Diracs = k-means**

*Application: Spectral clustering*
*for MNIST classification [Uw 2001]*

*Classif. Perf.*

- Twice faster than k-means
- 4 orders of magnitude more memory efficient

GMM

# Gaussian mixture models

# Gaussian mixture models



**GMM**

*d* = 10, **k = 20**

*Error*

*Size of database*

KL-div

- CL-OMPR
- EM1
- EM10

*Faster than EM (VLFeat's* gmm*)*

# Gaussian mixture models



**GMM**

d = 10, **k = 20**

*Error*

KL-div

*Faster than EM (VLFeat's gmm)*

*Size of database*

Legend: CL-OMPR, EM1, EM10

*Application: **speaker verification** [Reynolds 2000] (d=12, k=64)*
- EM on 300 000 vectors : **29.53**
- **20kB** sketch computed on **50GB** database: **28.96**

**Q: Theoretical guarantees ?**

- Inspired by Compressive Sensing:

  - 1: with the Restricted Isometry Property (RIP)

  - 2: with dual certificates

# Outline

**1** Information-preservation guarantees:
a RIP analysis

**2** Total variation regularization:
a dual certificate analysis

**3** Conclusion, outlooks

# Outline

① Information-preservation guarantees:
a RIP analysis
Joint work with **R. Gribonval, G. Blanchard, Y. Traonmilin**

② Total variation regularization:
a dual certificate analysis

③ Conclusion, outlooks

$$\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)$$

$\Phi$

$\mathcal{A}$

$+\hat{\mathbf{e}}$

$\hat{\mathbf{z}}$

$\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)$

$\pi^{\star}$

True distribution: $\boxed{x_1, ..., x_n \overset{i.i.d.}{\sim} \pi^{\star}}$

Sketch: $\boxed{\hat{\mathbf{z}} = \mathcal{A}\pi^{\star} + \hat{\mathbf{e}}}$

$$\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)$$

True distribution: $x_1, ..., x_n \overset{i.i.d.}{\sim} \pi^\star$

Sketch:

$$\hat{\mathbf{z}} = \mathcal{A}\pi^\star + \hat{\mathbf{e}}$$

- Estimation problem = **linear inverse problem** on measures

- **Extremely ill-posed !**

$$\hat{\mathbf{z}} = \hat{\mathbb{E}}\Phi(X)$$

True distribution: $x_1, ..., x_n \overset{i.i.d.}{\sim} \pi^\star$

Sketch: $\hat{\mathbf{z}} = \mathcal{A}\pi^\star + \hat{\mathbf{e}}$

- Estimation problem = **linear inverse problem** on measures

- **Extremely ill-posed !**

- *Feasibility? (information-preservation)*

$\mathcal{P}$

$\pi^{\star}$

$\mathscr{P}$

$\mathscr{S}$

$\pi^\star$

$\mathscr{S}$ : *Model set of « simple » distributions (eg. GMMs)*

$\mathcal{P}$

$\mathfrak{S}$

$\pi^{\star}$

$\boxed{\mathfrak{S} \text{ : Model set of « simple » distributions (eg. GMMs)}}$

$\mathcal{A}$

$\mathbb{C}^m$

$\mathcal{P}$

$\mathfrak{S}$

$\boxed{\mathfrak{S} : \text{Model set of « simple » distributions (eg. GMMs)}}$

$\pi^{\star}$

$\mathcal{A}$

$\mathbb{C}^m$

$\hat{\mathbf{z}}$ $+\hat{\mathbf{e}}$

$\mathcal{P}$

$\mathfrak{S}$

$\mathfrak{S}$ : *Model set of « simple » distributions (eg. GMMs)*

$\tilde{\pi}$

$\pi^{\star}$

$\Delta$

$\mathcal{A}$

$\mathbb{C}^m$

$\hat{\mathbf{z}}$ $+\hat{\mathbf{e}}$

## Goal

**Prove the existence of a *decoder* $\triangle$ robust to noise and stable to modeling error.**

*« Instance-optimal » decoder*

$\mathfrak{S}$ : *Model set of « simple »*
*distributions (eg. GMMs)*

**Goal**

Prove the existence of a *decoder* $\triangle$ robust
to **noise** and stable to **modeling error**.

*« Instance-optimal » decoder*

**Lower Restricted Isometry Property**

$$\|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2$$

$\mathfrak{S}$ : *Model set of « simple » distributions (eg. GMMs)*

$$\Delta(\hat{\mathbf{z}}) \in \arg\min_{\sigma \in \mathfrak{S}} \|\hat{\mathbf{z}} - \mathcal{A}\sigma\|_2$$

*Non-convex **generalized moment matching***

## Goal

**Prove the existence of a *decoder* $\Delta$ robust to noise and stable to modeling error.**

*« Instance-optimal » decoder*

## Lower Restricted Isometry Property

$$\|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2$$

$\mathcal{S}$ : *Model set of « simple » distributions (eg. GMMs)*

$$\Delta(\hat{\mathbf{z}}) \in \arg\min_{\sigma \in \mathcal{S}} \|\hat{\mathbf{z}} - \mathcal{A}\sigma\|_2$$

*Non-convex **generalized moment matching***

## Goal

**Prove the existence of a *decoder* $\Delta$ robust to noise and stable to modeling error.**

*« Instance-optimal » decoder*

## Lower Restricted Isometry Property

$$\|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2$$

**New goal: find/construct models $\mathcal{S}$ and operators $\mathcal{A}$ that satisfy the LRIP (w.h.p.)**

**Goal: LRIP** $\mathrm{w.h.p.}$ on $\mathcal{A}$, $\forall \sigma, \sigma' \in \mathfrak{S}$, $\|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2$.

# Appropriate metric

**Goal:** **LRIP** $\text{w.h.p. on } \mathcal{A}, \forall \sigma, \sigma' \in \mathfrak{S}, \|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2.$

Reproducing kernel:

$$\kappa(x, x') \quad \left\langle \; \rule{0.3em}{2em} \; , \; \rule{0.3em}{2em} \; \right\rangle$$

**Goal: LRIP** $\mathrm{w.h.p.\ on\ }\mathcal{A},\ \forall \sigma, \sigma' \in \mathfrak{S},\ \|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2.$



Reproducing kernel:

$\kappa(x, x')$

**Kernel mean**

$\kappa(\pi, \pi') = \mathbb{E}\kappa(X, X')$

**Goal: LRIP** $\mathrm{w.h.p.}$ on $\mathcal{A}$, $\forall \sigma, \sigma' \in \mathfrak{S}$, $\|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2$.

Reproducing kernel:

$\kappa(x, x')$ $\left\langle \phantom{|} , \phantom{|} \right\rangle$



$\Phi$ : random features *[Rahimi2007]* to approximate $\kappa$

**Kernel mean**

$\kappa(\pi, \pi') = \mathbb{E}\kappa(X, X')$

**Goal: LRIP** $\mathrm{w.h.p.}$ $\mathrm{on}$ $\mathcal{A}$, $\forall \sigma, \sigma' \in \mathfrak{S}$, $\|\sigma - \sigma'\| \lesssim \|\mathcal{A}\sigma - \mathcal{A}\sigma'\|_2.$

Reproducing kernel:

$\kappa(x, x')$ $\left\langle \quad , \quad \right\rangle$

$\Phi$ : random features *[Rahimi2007]* to approximate $\kappa$

$\mathcal{A}\pi = \mathbb{E}_\pi \Phi(X)$

**Kernel mean**

$\kappa(\pi, \pi') = \mathbb{E}\kappa(X, X')$

$\left\langle \quad , \quad \right\rangle$

**Basis for LRIP**

$$\|\pi - \pi'\|_\kappa^2 \approx \|\mathcal{A}\pi - \mathcal{A}\pi'\|_2^2$$

**Reformulation of the LRIP**

**Goal: LRIP** $\quad \|\sigma - \sigma'\|_\kappa \lesssim \|\mathcal{A}(\sigma - \sigma')\|_2$

**Reformulation of the LRIP**

**Goal: LRIP** $\quad \|\sigma - \sigma'\|_\kappa \lesssim \|\mathcal{A}(\sigma - \sigma')\|_2$

$$\Leftrightarrow 1 \lesssim \left\|\mathcal{A}\left(\frac{\sigma - \sigma'}{\|\sigma - \sigma'\|_\kappa}\right)\right\|_2$$

**Reformulation of the LRIP**

**Goal: LRIP** $\quad \|\sigma - \sigma'\|_\kappa \lesssim \|\mathcal{A}(\sigma - \sigma')\|_2$

$$\Leftrightarrow \boxed{1 \lesssim \left\|\mathcal{A}\left(\frac{\sigma - \sigma'}{\|\sigma - \sigma'\|_\kappa}\right)\right\|_2}$$

**Definition: Normalized Secant set**

$$\mathcal{S}_\mathfrak{S} = \left\{\frac{\sigma - \sigma'}{\|\sigma - \sigma'\|_\kappa}; \ \sigma, \sigma' \in \mathfrak{S}\right\}$$

$\mathcal{M}$

$\mathcal{S}_\mathfrak{S}$

**Reformulation of the LRIP**

Goal: LRIP $\quad \|\sigma - \sigma'\|_\kappa \lesssim \|\mathcal{A}(\sigma - \sigma')\|_2$

$$\Leftrightarrow \boxed{1 \lesssim \|\mathcal{A}(\frac{\sigma - \sigma'}{\|\sigma - \sigma'\|_\kappa})\|_2}$$

**Definition: Normalized Secant set**

$$\mathcal{S}_{\mathfrak{S}} = \left\{ \frac{\sigma - \sigma'}{\|\sigma - \sigma'\|_\kappa}; \ \sigma, \sigma' \in \mathfrak{S} \right\}$$

**New goal**

With high probability on $\mathcal{A}$ :

*for all* $s \in \mathcal{S}_{\mathfrak{S}}, 1 \lesssim \|\mathcal{A}s\|_2$ .



$\mathcal{M}$

$\mathcal{S}_{\mathfrak{S}}$

$\mathcal{A}$

$\mathbb{C}^m$

0

**Goal: LRIP** $\quad$ w.h.p. on $\mathcal{A}$, $\forall s \in \mathcal{S}_{\mathfrak{S}}$, $1 \lesssim \|\mathcal{A}s\|_2$.

**Goal: LRIP** $\quad$ w.h.p. on $\mathcal{A}$, $\forall s \in \mathcal{S}_{\mathfrak{S}}$, $1 \lesssim \|\mathcal{A}s\|_2$.

① **Pointwise LRIP: Concentration inequality** $\quad$ $\forall s$, w.h.p. on $\mathcal{A}$, LRIP.

**Goal: LRIP** $\quad$ w.h.p. on $\mathcal{A}$, $\forall s \in \mathcal{S}_{\mathfrak{S}}$, $1 \lesssim \|\mathcal{A}s\|_2$.

① **Pointwise LRIP: Concentration inequality** $\qquad$ $\forall s$, w.h.p. on $\mathcal{A}$, LRIP.

② **Extension to LRIP: covering numbers**



w.h.p. on $\mathcal{A}$, $\forall s$, LRIP.

**Main hypothesis**

The *normalized secant set* $\mathcal{S}(\mathfrak{S})$ has finite covering numbers.

# Main result

**Main hypothesis**

The *normalized secant set* $\mathcal{S}(\mathfrak{S})$ has finite covering numbers.

**Result**

For $\boxed{m \geq C \ \times \ \log(\text{cov. num.})}$ ,

**Main hypothesis**

The *normalized secant set* $\mathcal{S}(\mathfrak{S})$ has finite covering numbers.

**Result**

For $\boxed{m \geq C \ \times \ \log(\text{cov. num.})}$ ,

Quality of pointwise LRIP

Dimensionality of the model

**Main hypothesis**

The *normalized secant set* $\mathcal{S}(\mathfrak{S})$ has finite covering numbers.

**Result**

For $\boxed{m \geq C \ \times \ \log(\text{cov. num.})}$ ,

Quality of pointwise LRIP      Dimensionality of the model

W.h.p.

$$\|\pi^\star - \Delta(\hat{\mathbf{z}})\| \leq d(\pi^\star, \mathfrak{S}) + \mathcal{O}(1/\sqrt{n})$$

**Main hypothesis**

The *normalized secant set* $\mathcal{S}(\mathfrak{S})$ has finite covering numbers.

**Result**

For $\boxed{m \geq C \times \log(\text{cov. num.})}$ ,

Quality of pointwise LRIP

Dimensionality of the model

W.h.p.

Modeling error

Empirical noise

$$\|\pi^\star - \Delta(\hat{\mathbf{z}})\| \leq d(\pi^\star, \mathfrak{S}) + \mathcal{O}(1/\sqrt{n})$$

**Main hypothesis**

The *normalized secant set* $\mathcal{S}(\mathfrak{S})$ has finite covering numbers.

**Result**

For $\boxed{m \geq C \times \log(\text{cov. num.})}$ ,

Quality of pointwise LRIP    Dimensionality of the model

W.h.p.

*Modeling error*    *Empirical noise*

$$\|\pi^\star - \Delta(\hat{\mathbf{z}})\| \leq d(\pi^\star, \mathfrak{S}) + \mathcal{O}(1/\sqrt{n})$$

- **Classic Compressive Sensing**: finite dimension: **Known**
- **Here**: infinite dimension: **Technical**

**k-means with mixtures of Diracs**

**k-means with mixtures of Diracs**

**Hypotheses**
- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

**k-means with mixtures of Diracs**

**Hypotheses**

*(no assumption on the **data**)*

- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

## k-means with mixtures of Diracs

### Hypotheses
*(no assumption on the **data**)*

- $\varepsilon$ - separated centroids
- $M$- bounded domain for centroids

### Sketch
- *Adjusted* Random Fourier features *(for technical reasons)*

## k-means with mixtures of Diracs

### Hypotheses

*(no assumption on the **data**)*

- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

### Sketch

- *Adjusted* Random Fourier features *(for technical reasons)*

### Result

- W.r.t. **k-means usual cost (SSE)**

## k-means with mixtures of Diracs

### Hypotheses

*(no assumption on the **data**)*

- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

### Sketch

- *Adjusted* Random Fourier features *(for technical reasons)*

### Result

- W.r.t. **k-means usual cost (SSE)**

### Sketch size

$$m \geq \mathcal{O}\left(k^2 d \cdot \texttt{polylog}(k, d) \log(M/\varepsilon)\right)$$

# Application

## k-means with mixtures of Diracs

### Hypotheses
*(no assumption on the **data**)*

- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

### Sketch
- *Adjusted* Random Fourier features *(for technical reasons)*

### Result
- W.r.t. **k-means usual cost (SSE)**

### Sketch size

$$m \geq \mathcal{O}\left(k^2 d \cdot \mathrm{polylog}(k, d) \log(M/\varepsilon)\right)$$

## GMM with known covariance

# Application

## k-means with mixtures of Diracs

### Hypotheses
*(no assumption on the **data**)*
- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

### Sketch
- *Adjusted* Random Fourier features *(for technical reasons)*

### Result
- W.r.t. **k-means usual cost (SSE)**

### Sketch size

$$m \geq \mathcal{O}\left(k^2 d \cdot \text{polylog}(k, d) \log(M/\varepsilon)\right)$$

## GMM with known covariance

### Hypotheses
- Sufficiently separated means
- Bounded domain for means

# Application

<table>
<tr><td>

**k-means with mixtures of Diracs**

**Hypotheses** *(no assumption on the **data**)*
- $\varepsilon$ - separated centroids
- $M$- bounded domain for centroids

**Sketch**
- *Adjusted* Random Fourier features *(for technical reasons)*

**Result**
- W.r.t. **k-means usual cost (SSE)**

**Sketch size**

$$m \geq \mathcal{O}\left(k^2 d \cdot \texttt{polylog}(k, d) \log(M/\varepsilon)\right)$$

</td><td>

**GMM with known covariance**

**Hypotheses**
- Sufficiently separated means
- Bounded domain for means

**Sketch**
- Fourier features

</td></tr>
</table>

CFM INSIGHT.DATA.CLARITY.  ENS  PSL★ RESEARCH UNIVERSITY PARIS

# Application

<table>
<tr><th>k-means with mixtures of Diracs</th><th>GMM with known covariance</th></tr>
</table>

**Hypotheses**  *(no assumption on the **data**)*
- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

**Sketch**
- *Adjusted* Random Fourier features *(for technical reasons)*

**Result**
- W.r.t. **k-means usual cost (SSE)**

**Sketch size**

$$m \geq \mathcal{O}\left(k^2 d \cdot \texttt{polylog}(k, d) \log(M/\varepsilon)\right)$$

**Hypotheses**
- Sufficiently separated means
- Bounded domain for means

**Sketch**
- Fourier features

**Result**
- With respect to **log-likelihood**

# Application

| k-means with mixtures of Diracs | GMM with known covariance |
|---|---|

### k-means with mixtures of Diracs

**Hypotheses** *(no assumption on the **data**)*
- $\varepsilon$ - separated centroids
- $M$ - bounded domain for centroids

**Sketch**
- *Adjusted* Random Fourier features *(for technical reasons)*

**Result**
- W.r.t. **k-means usual cost (SSE)**

**Sketch size**

$$m \geq \mathcal{O}\left(k^2 d \cdot \mathrm{polylog}(k, d) \log(M/\varepsilon)\right)$$

### GMM with known covariance

**Hypotheses**
- Sufficiently separated means
- Bounded domain for means

**Sketch**
- Fourier features

**Result**
- With respect to **log-likelihood**

**Sketch size**

$$m \geq \mathcal{O}(k^2 d \cdot \mathrm{polylog}(k, d))$$

**With the RIP analysis:**

# Summary

With the RIP analysis:

- **Moment matching**: best decoder possible (instance optimal)
  - Information-preservation guarantees

# Summary

**With the RIP analysis:**

- **Moment matching**: best decoder possible (instance optimal)
  - Information-preservation guarantees

- Fine control on modeling error, noise, and metrics
  - **Can incorporate k-means cost or log-likelihood**

# Summary

**With the RIP analysis:**

- **Moment matching**: best decoder possible (instance optimal)
  - Information-preservation guarantees

- Fine control on modeling error, noise, and metrics
  - **Can incorporate k-means cost or log-likelihood**

**Compressive Sensing:**

**With the RIP analysis:**

- **Moment matching**: best decoder possible (instance optimal)
  - Information-preservation guarantees

- Fine control on modeling error, noise, and metrics
  - **Can incorporate k-means cost or log-likelihood**

**Compressive Sensing:**

- Random, dimensionality-reducing operator ✓

# Summary

**With the RIP analysis:**

- **Moment matching**: best decoder possible (instance optimal)
    - Information-preservation guarantees

- Fine control on modeling error, noise, and metrics
    - **Can incorporate k-means cost or log-likelihood**

**Compressive Sensing:**

- Random, dimensionality-reducing operator ✓

- Sparsity ✓

# Summary

## With the RIP analysis:

- **Moment matching**: best decoder possible (instance optimal)
  - Information-preservation guarantees

- Fine control on modeling error, noise, and metrics
  - **Can incorporate k-means cost or log-likelihood**

## Compressive Sensing:

- Random, dimensionality-reducing operator ✓

- Sparsity ✓

- The information is preserved ✓

# Summary

---

**With the RIP analysis:**

- **Moment matching**: best decoder possible (instance optimal)
    - Information-preservation guarantees

- Fine control on modeling error, noise, and metrics
    - **Can incorporate k-means cost or log-likelihood**

---

**Compressive Sensing:**

- Random, dimensionality-reducing operator ✓

- Sparsity ✓

- The information is preserved ✓

- **Convex relaxation?** ✗

---

# Outline

① Information-preservation guarantees:
a RIP analysis

② Total variation regularization:
a dual certificate analysis
Joint work with **C. Poon, G. Peyré**

③ Conclusion, outlooks

Nicolas Keriven

**Previously: RIP analysis**

*Minimization: moment matching*

$$\min_\theta \|\mathcal{A}(\sum w_i \pi_{\theta_i}) - \hat{\mathbf{z}}\|_2$$

# Total Variation regularization

**Previously: RIP analysis**

*Minimization: moment matching*

$$\min_\theta \left\| \mathcal{A}\left(\sum w_i \pi_{\theta_i}\right) - \hat{\mathbf{z}} \right\|_2$$

- Must know $k$

- **Non-convex !**

# Total Variation regularization

**Previously: RIP analysis**

*Minimization: moment matching*

$$\min_\theta \left\| \mathcal{A}\left( \sum w_i \pi_{\theta_i} \right) - \hat{\mathbf{z}} \right\|_2$$

- Must know $k$

- **Non-convex !**

**Convex relaxation (« *super resolution* »)**

$$\min_\mu \frac{1}{2} \| \Psi\mu - \hat{\mathbf{z}} \|_2 + \lambda \| \mu \|_{\mathrm{TV}}$$

- $\mu$ : Radon measure

- $\Psi\mu = \int (\mathcal{A}\pi_\theta) d\mu(\theta)$

- $\| \cdot \|_{\mathrm{TV}}$ : Total variation (« L1 norm »)

# Total Variation regularization

**Previously: RIP analysis**

*Minimization: moment matching*

$$\min_\theta \|\mathcal{A}(\sum w_i \pi_{\theta_i}) - \hat{\mathbf{z}}\|_2$$

- Must know $k$

- **Non-convex !**

**Convex relaxation (« *super resolution* »)**

$$\min_\mu \frac{1}{2}\|\Psi\mu - \hat{\mathbf{z}}\|_2 + \lambda\|\mu\|_{\mathrm{TV}}$$

- $\mu$ : Radon measure

- $\Psi\mu = \int(\mathcal{A}\pi_\theta)d\mu(\theta)$

- $\|\cdot\|_{\mathrm{TV}}$ : Total variation (« L1 norm »)

**Convex**:
- can be handled by eg Frank-Wolfe algorithm *[Boyd 2015]*, or in some cases as a SDP

# Total Variation regularization

## Previously: RIP analysis

*Minimization: moment matching*

$$\min_\theta \|\mathcal{A}(\sum w_i \pi_{\theta_i}) - \hat{\mathbf{z}}\|_2$$

- Must know $k$

- **Non-convex !**

## Convex relaxation (« *super resolution* »)

$$\min_\mu \frac{1}{2}\|\Psi\mu - \hat{\mathbf{z}}\|_2 + \lambda\|\mu\|_{\mathrm{TV}}$$

- $\mu$ : Radon measure

- $\Psi\mu = \int (\mathcal{A}\pi_\theta)d\mu(\theta)$

- $\|\cdot\|_{\mathrm{TV}}$ : Total variation (« L1 norm »)

**Convex**:
- can be handled by eg Frank-Wolfe algorithm *[Boyd 2015]*, or in some cases as a SDP

**Questions**:
- Is the measure $\mu$ sparse ? $\mu = \sum \tilde{w}_i \delta_{\tilde{\theta}_i}$

- Does it have the right number of components ?

- Does it recover the true $w_i, \theta_i$ ?

# A bit of convex analysis

**Intuition**: first order conditions: $\mu_0$ solution $\quad \Longleftrightarrow \quad \frac{1}{\lambda}\Psi^{\star}(\Psi\mu_0 - \hat{\mathbf{z}}) \in \partial\|\mu_0\|_{\mathrm{TV}}$

# A bit of convex analysis

**Intuition**: first order conditions: $\mu_0$ solution $\qquad \Longleftrightarrow \qquad \frac{1}{\lambda}\Psi^\star(\Psi\mu_0 - \hat{\mathbf{z}}) \in \partial\|\mu_0\|_{\mathrm{TV}}$

**Def. : Dual certificate** ( = Lagrange multiplier in the noiseless case…)

$$\eta \in \mathrm{Im}(\Psi^\star) \cap \partial\|\mu_0\|_{\mathrm{TV}}$$

# A bit of convex analysis

**Intuition**: first order conditions: $\mu_0$ solution $\iff$ $\frac{1}{\lambda}\Psi^\star(\Psi\mu_0 - \hat{\mathbf{z}}) \in \partial\|\mu_0\|_{\mathrm{TV}}$

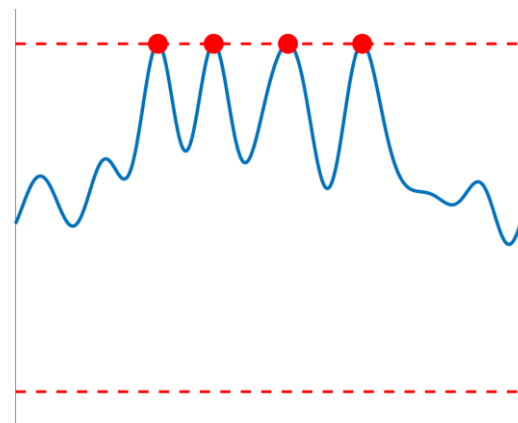**Def. : Dual certificate** ( = Lagrange multiplier in the noiseless case...)

$$\eta \in \mathrm{Im}(\Psi^\star) \cap \partial\|\mu_0\|_{\mathrm{TV}}$$

**What is a dual certificate?**

# A bit of convex analysis

**Intuition**: first order conditions: $\mu_0$ solution $\iff$ $\frac{1}{\lambda}\Psi^\star(\Psi\mu_0 - \hat{\mathbf{z}}) \in \partial\|\mu_0\|_{\mathrm{TV}}$

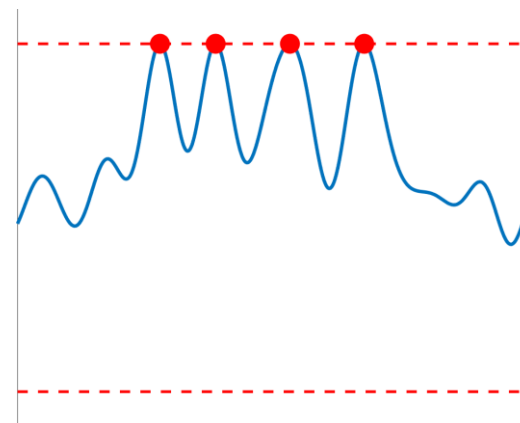**Def. : Dual certificate** ( = Lagrange multiplier in the noiseless case…)

$$\eta \in \mathrm{Im}(\Psi^\star) \cap \partial\|\mu_0\|_{\mathrm{TV}}$$

**What is a dual certificate?**

$$\eta(\theta) = \langle \mathbf{h}, \mathcal{A}\pi_\theta \rangle$$

Such that:

- $\eta(\theta_i) = 1$
- $|\eta(\theta)| < 1$  otherwise
- $\nabla^2\eta(\theta_i) \prec 0$

# A bit of convex analysis

**Intuition**: first order conditions: $\mu_0$ solution $\iff$ $\frac{1}{\lambda}\Psi^\star(\Psi\mu_0 - \hat{\mathbf{z}}) \in \partial\|\mu_0\|_{\mathrm{TV}}$

**Def. : Dual certificate** ( = Lagrange multiplier in the noiseless case…)

$$\eta \in \mathrm{Im}(\Psi^\star) \cap \partial\|\mu_0\|_{\mathrm{TV}}$$

**What is a dual certificate?**

$$\eta(\theta) = \langle \mathbf{h}, \mathcal{A}\pi_\theta \rangle$$

Such that:

- $\eta(\theta_i) = 1$
- $|\eta(\theta)| < 1$ otherwise
- $\nabla^2\eta(\theta_i) \prec 0$

*Ensures uniqueness and robustness…*

# Strategy: going back to random features

**Step 1: study full kernel**

**Step 1: study full kernel**

$$\bar{\eta} \in \mathrm{Span}\left\{\kappa(\theta_i, \cdot), \partial_1 \kappa(\theta_i, \cdot)\right\} \subset \mathrm{Im}(\mathbb{E}\Psi^\star)$$

**Step 1: study full kernel**

$$\bar{\eta} \in \mathrm{Span}\left\{\kappa(\theta_i, \cdot), \partial_1 \kappa(\theta_i, \cdot)\right\} \subset \mathrm{Im}(\mathbb{E}\Psi^\star)$$

**Assumptions**:
- Kernel « well-behaved »
- $\theta_i$ sufficiently separated



$m = \infty$

## Step 1: study full kernel

$$\bar{\eta} \in \mathrm{Span}\left\{\kappa(\theta_i, \cdot), \partial_1 \kappa(\theta_i, \cdot)\right\} \subset \mathrm{Im}(\mathbb{E}\Psi^\star)$$

**Assumptions**:
- Kernel « well-behaved »
- $\theta_i$ sufficiently separated


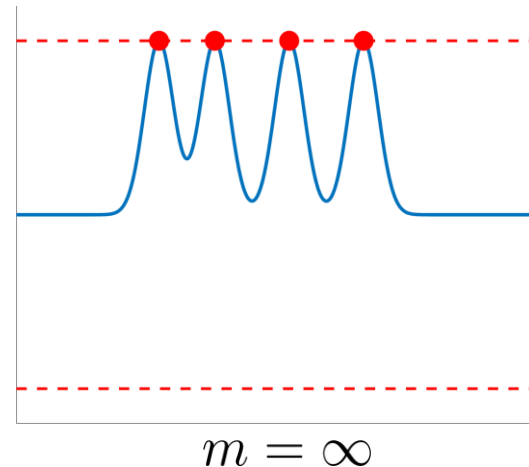
$m = \infty$
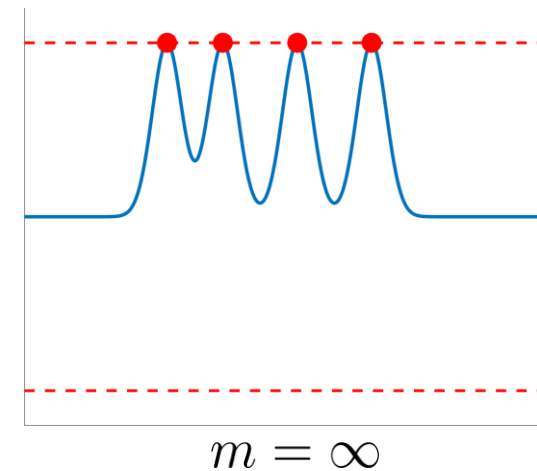
## Step 2: bounding the deviations

# Strategy: going back to random features

## Step 1: study full kernel

$$\bar{\eta} \in \operatorname{Span}\left\{\kappa(\theta_i, \cdot), \partial_1 \kappa(\theta_i, \cdot)\right\} \subset \operatorname{Im}(\mathbb{E}\Psi^\star)$$

**Assumptions**:
- Kernel « well-behaved »
- $\theta_i$ sufficiently separated



$$m = \infty$$

## Step 2: bounding the deviations

- Pointwise deviation (concentration ineq.)
- Covering numbers

# Strategy: going back to random features

## Step 1: study full kernel

$$\bar{\eta} \in \mathrm{Span}\left\{\kappa(\theta_i, \cdot), \partial_1 \kappa(\theta_i, \cdot)\right\} \subset \mathrm{Im}(\mathbb{E}\Psi^\star)$$
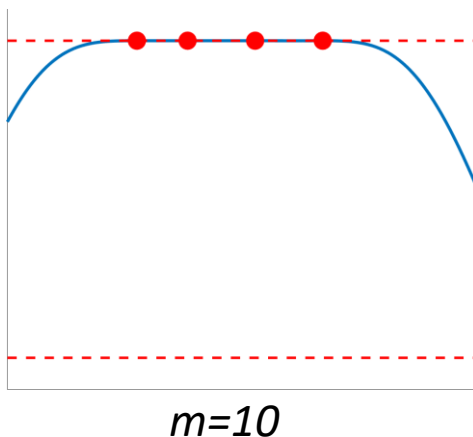
**Assumptions**:
- Kernel « well-behaved »
- $\theta_i$ sufficiently separated



$m = \infty$

## Step 2: bounding the deviations

- Pointwise deviation (concentration ineq.)
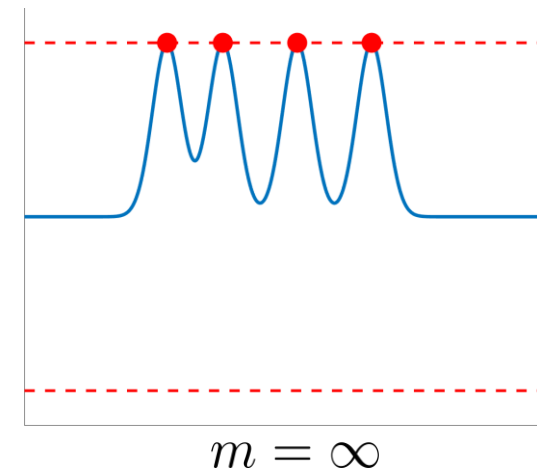- Covering numbers



*m=10*

# Strategy: going back to random features

## Step 1: study full kernel

$$\bar{\eta} \in \mathrm{Span}\left\{\kappa(\theta_i, \cdot), \partial_1 \kappa(\theta_i, \cdot)\right\} \subset \mathrm{Im}(\mathbb{E}\Psi^\star)$$
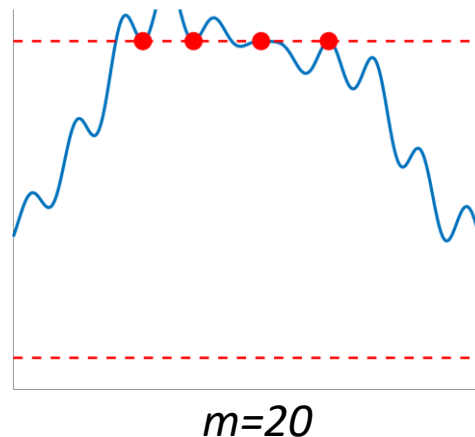
**Assumptions**:
- Kernel « well-behaved »
- $\theta_i$ sufficiently separated

$m = \infty$

## Step 2: bounding the deviations

- Pointwise deviation (concentration ineq.)
- Covering numbers

*m=10*

*m=20*

## Step 1: study full kernel

$$\bar{\eta} \in \mathrm{Span}\left\{\kappa(\theta_i, \cdot), \partial_1 \kappa(\theta_i, \cdot)\right\} \subset \mathrm{Im}(\mathbb{E}\Psi^\star)$$
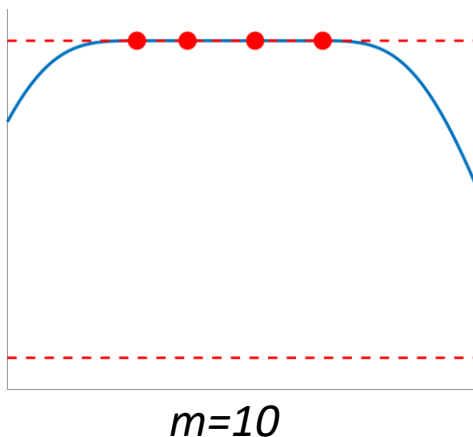
**Assumptions**:
- Kernel « well-behaved »
- $\theta_i$ sufficiently separated

$m = \infty$

## Step 2: bounding the deviations

- Pointwise deviation (concentration ineq.)
- Covering numbers

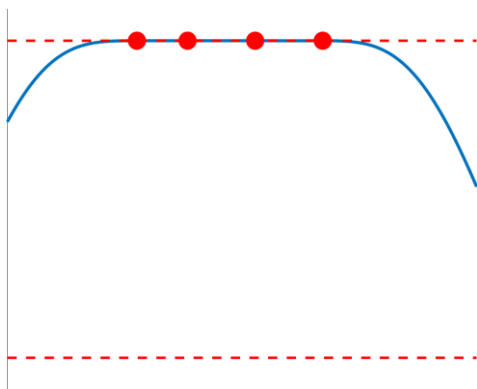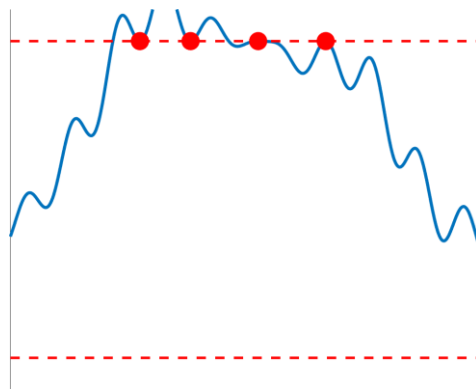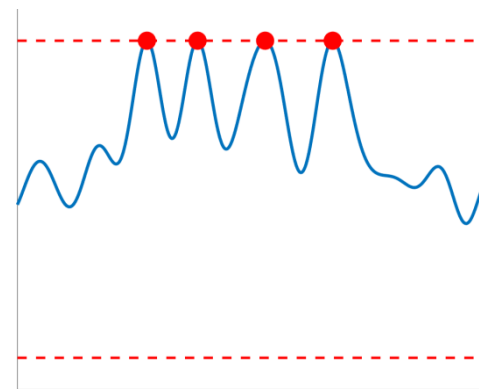*m=10*  *m=20*  *m=50*

Assumption: data are *actually* drawn from a GMM...

**1: Ideal scaling in sparsity**

Assumption: data are *actually* drawn from a GMM…

## 1: Ideal scaling in sparsity

$$m \geq \mathcal{O}(kd^4 \cdot \texttt{polylog}(k, d))$$

Assumption: data are *actually* drawn from a GMM...

## 1: Ideal scaling in sparsity

$$m \geq \mathcal{O}(k d^4 \cdot \mathrm{polylog}(k, d))$$

↑

*In progress...*

Assumption: data are *actually* drawn from a GMM…

## 1: Ideal scaling in sparsity

$$m \geq \mathcal{O}(k d^4 \cdot \text{polylog}(k, d))$$

↑
*In progress…*

- $\tilde{\mu}$ **not necessarily sparse**, but:

- Mass of $\tilde{\mu}$ concentrated around true $\theta_i$

- *Proof*: infinite-dimensional golfing scheme (new)

Assumption: data are *actually* drawn from a GMM...

## 1: Ideal scaling in sparsity

$$m \geq \mathcal{O}(kd^4 \cdot \texttt{polylog}(k,d))$$

↑

*In progress...*

- $\tilde{\mu}$ **not necessarily sparse**, but:

- Mass of $\tilde{\mu}$ concentrated around true $\theta_i$

- *Proof*: infinite-dimensional golfing scheme (new)

## 2: *Minimal norm* certificate

[Duval, Peyré 2015]

$$m \geq \mathcal{O}(k^2 d^3 \cdot \texttt{polylog}(k,d))$$

↑

*In progress...*

# Results for separated GMM

Assumption: data are *actually* drawn from a GMM…

## 1: Ideal scaling in sparsity

$$m \geq \mathcal{O}(k d^4 \cdot \texttt{polylog}(k, d))$$

↑
*In progress…*

- $\tilde{\mu}$ ***not necessarily sparse***, but:

- Mass of $\tilde{\mu}$ concentrated around true $\theta_i$

- *Proof*: infinite-dimensional golfing scheme (new)

## 2: *Minimal norm* certificate

*[Duval, Peyré 2015]*

$$m \geq \mathcal{O}(k^2 d^3 \cdot \texttt{polylog}(k, d))$$

↑
*In progress…*

- when *n* high enough: $\tilde{\mu}$ **sparse, with right number of components**

- $\tilde{\theta}_i \xrightarrow[n \to \infty]{} \theta_i$

- Proof: adaptation of *[Tang, Recht 2013] (constructive!)*
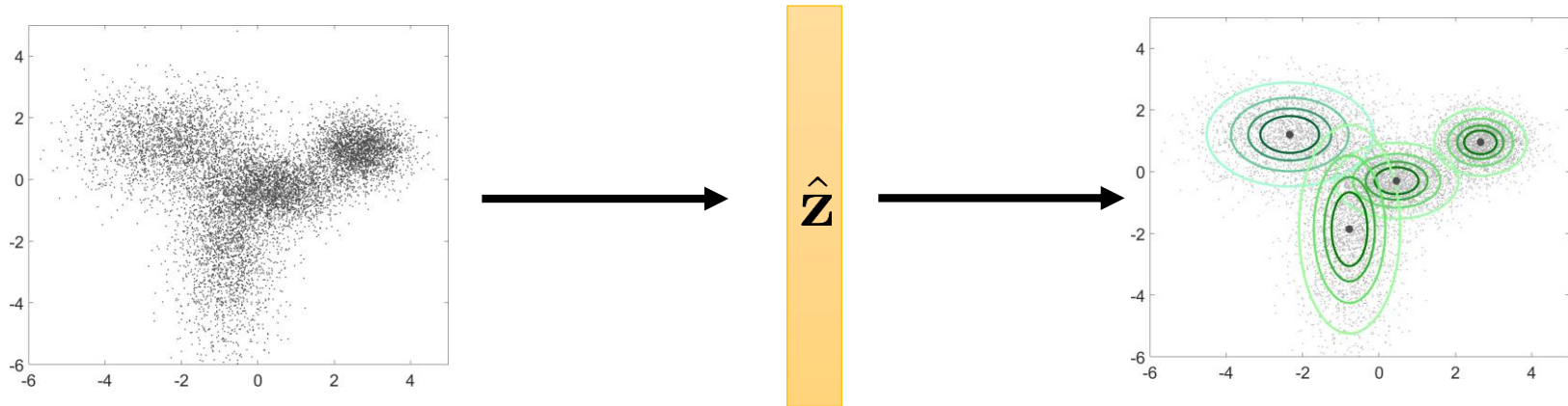
① Information-preservation guarantees:
a RIP analysis

② Total variation regularization:
a dual certificate analysis

③ Conclusion, outlooks

- Sketching :
  - Streaming, distributed learning

  - Original view on data compression and generalized moments

  - Combines random features and kernel mean with infinite dimensional Compressive sensing

- **RIP analysis**
  - Information preservation guarantees
  - Fine control on noise, modeling error (instance optimal decoder) and recovery metrics
  - Necessary and sufficient conditions
  - But: Non-convex minimization

# Summary, outlooks

- **RIP analysis**
  - Information preservation guarantees
  - Fine control on noise, modeling error (instance optimal decoder) and recovery metrics
  - Necessary and sufficient conditions
  - But: Non-convex minimization

- **Dual certificate analysis**
  - Convex minimization
  - Does not handle modelling error
  - In some cases, automatically guess the right number of components

# Summary, outlooks

- **RIP analysis**
  - Information preservation guarantees
  - Fine control on noise, modeling error (instance optimal decoder) and recovery metrics
  - Necessary and sufficient conditions
  - But: Non-convex minimization

- **Dual certificate analysis**
  - Convex minimization
  - Does not handle modelling error
  - In some cases, automatically guess the right number of components

- **Outlooks**
  - Algorithms for TV minimization
  - Other features $\Phi$ (not necessarily random...)
  - Other « sketched » learning tasks
  - Multilayer sketches ?

# Thank you !

- Keriven, Bourrier, Gribonval, Pérez. **Sketching for Large-Scale Learning of Mixture Models** *Information & Inference: a Journal of the IMA,* 2017. <arXiv:1606.02838>

- Keriven, Tremblay, Traonmilin, Gribonval. **Compressive k-means** *ICASSP,* 2017.

- Gribonval, Blanchard, Keriven, Traonmilin. **Compressive Statistical Learning with Random Feature Moments**. *Preprint* 2017. <arXiv:1706.07180>

- Keriven. **Sketching for Large-Scale Learning of Mixture Models**. *PhD Thesis.* <tel-01620815>

- Poon, Keriven, Peyré. **A Dual Certificates Analysis of Compressive Off-the-Grid Recovery.** *Submitted*

- **Code**: sketchml.gforge.inria.fr,
        github: nkeriven