

Inria



ENS DE LYON

Compressive learning and sketching for large-scale machine learning

Titouan Vayer

INRIA Lyon - Ockham Team

titouan.vayer@inria.fr

Motivations of this talk

- A modern fairy tale

Fashion Trend Forecasting with AI


JESSICA MAGALIT - NOVEMBER 9, 2021

🔗 0 💬 0

See What's Next

With T-Fashion's AI-powered trend forecasting platform, grasp trend dynamics through billions of interactions taking place online.

Receive customized fashion analytics and data-driven trend insights to produce/buy the right product at the right time.



Sunglasses

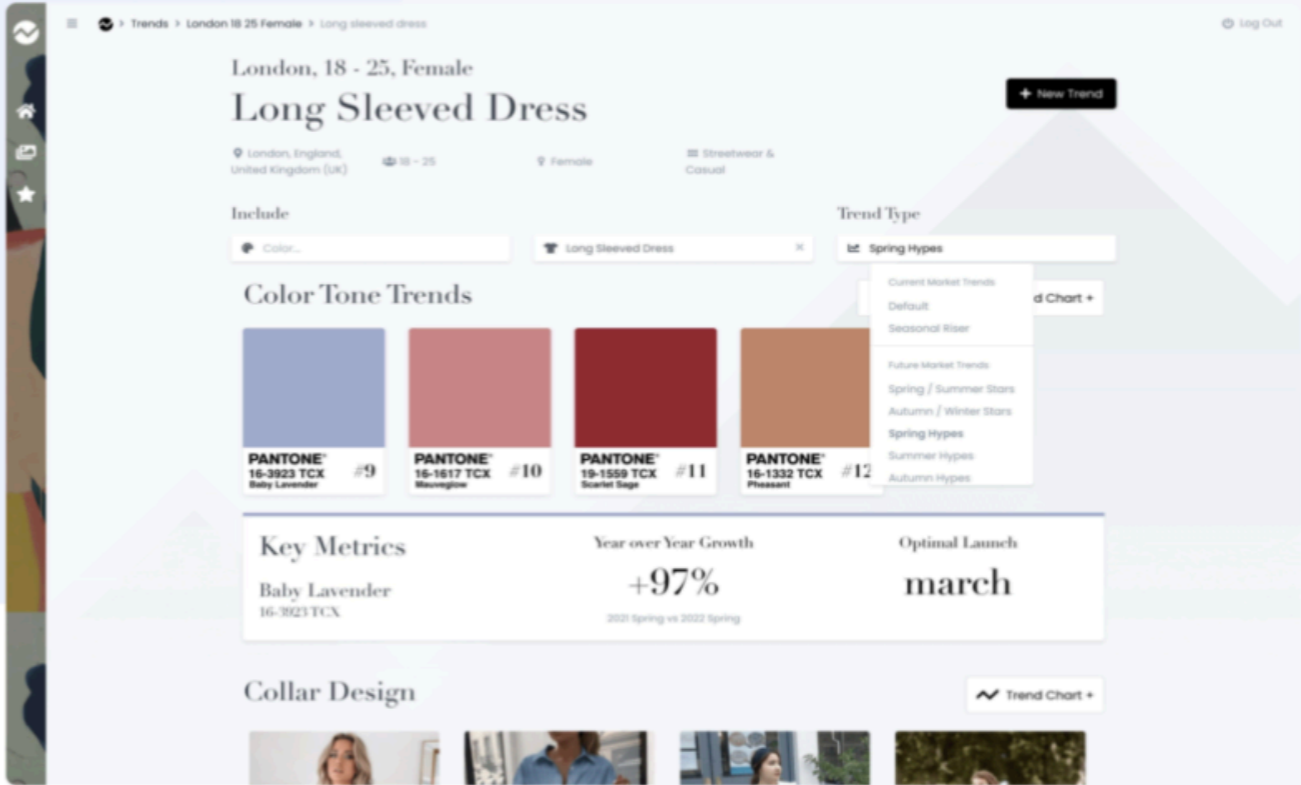
Gender: Female
Age: 29
Location: Kyiv, Ukraine

Outwear

- Long sleeved
- Notched collar
- Classic Blue

Dress

- Short sleeved
- Straight Neck
- Knee Length
- Lily White - Classic Blue



London, 18 - 25, Female
Long Sleeved Dress

London, England, United Kingdom (UK) | 18 - 25 | Female | Streetwear & Casual

Include: Color... | Long Sleeved Dress | Trend Type: Spring Hypes

Color Tone Trends

PANTONE	TCX	#
Baby Lavender	16-3923	#9
Mauveglow	16-1817	#10
Scarlet Sage	19-1559	#11
Pheasant	16-1332	#12

Key Metrics

- Baby Lavender 16-3923 TCX
- Year over Year Growth: **+97%** (2021 Spring vs 2022 Spring)
- Optimal Launch: **march**

Collar Design

Trend Chart +

| Motivations of this talk

- A modern fairy tale



Motivations of this talk

■ A modern fairy tale



Predicting the future isn't magic, it's artificial intelligence.

DAVE WATERS

HR Technology

An AI revolution is underway, predicts Kirthiga Reddy

Hiring without bias, bringing diversity into organisations, eliminating unconscious bias: AI is game-changing, says the president of Athena SPACs.

The artificial intelligence revolution



18th March 2019
Electronic Specifier
Lanna Deamer



| Motivations of this talk

- A modern fairy tale



Motivations of this talk

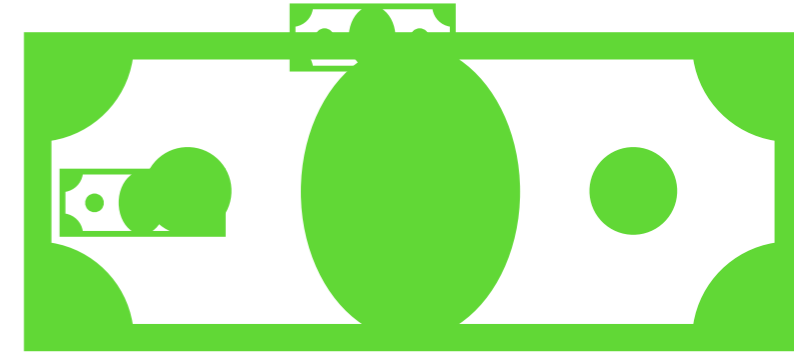
- A modern fairy tale



The internet
1e5 TB

Motivations of this talk

■ A modern fairy tale



predict



The internet
1e5 TB



| Motivations of this talk

- A modern fairy tale



| Motivations of this talk

- A modern fairy tale



| Motivations of this talk

- A modern fairy tale



The internet
1e5 TB

Motivations of this talk

■ A modern fairy tale

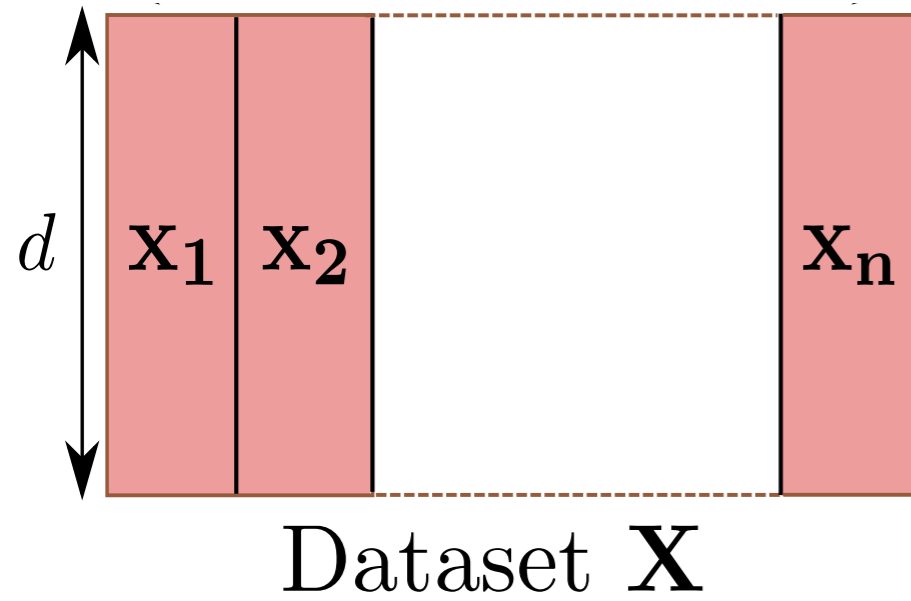


Statistical correlation ...



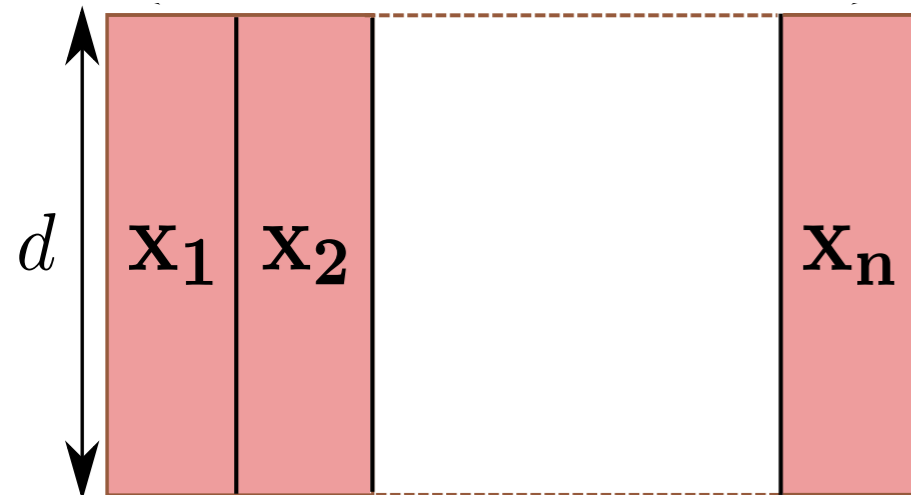
| Motivations of this talk

- **Context:** this will not be a fairy tale



| Motivations of this talk

- **Context:** this will not be a fairy tale



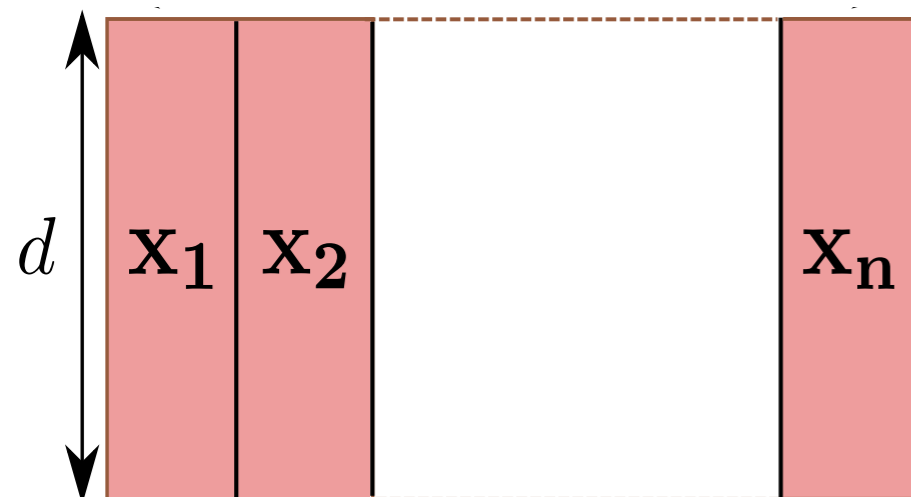
Dataset X

a sample (e.g. vector, image, embedding of words)



Motivations of this talk

■ **Context:** this will not be a fairy tale



Dataset X

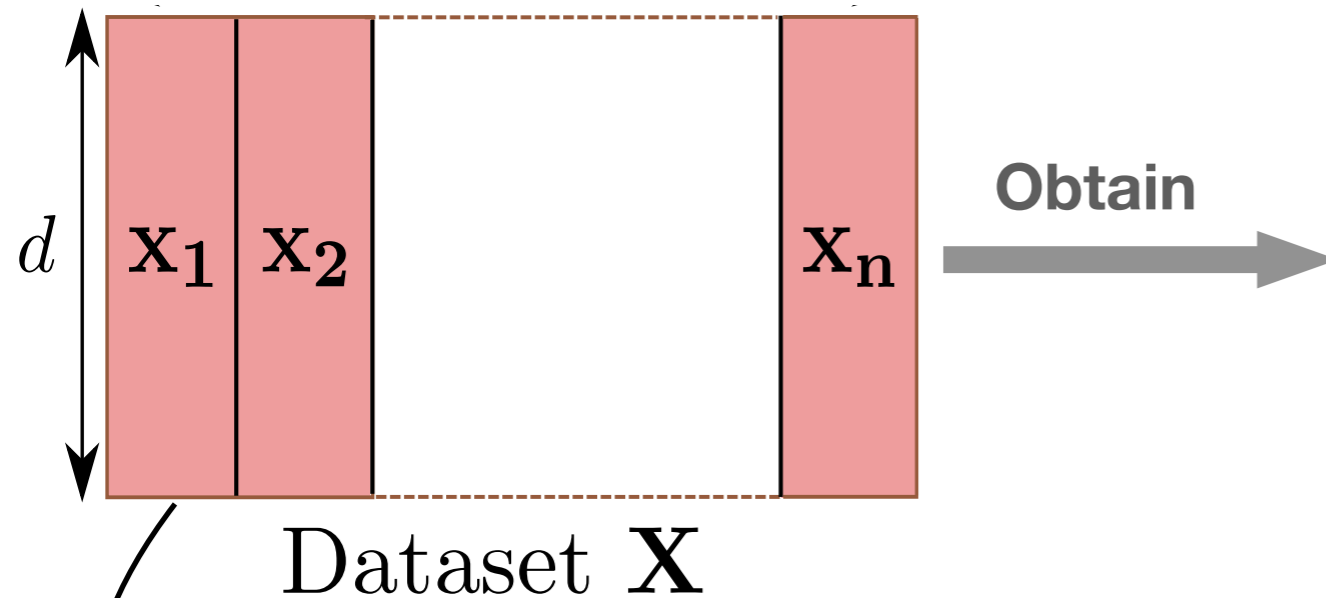
a sample (e.g. vector, image, embedding of words)



we are not like John

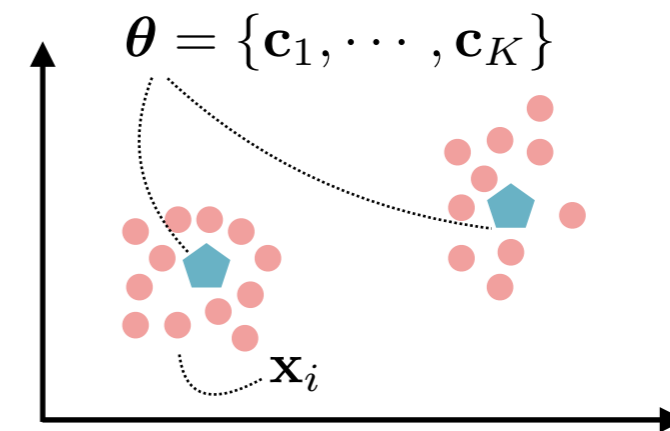
Motivations of this talk

Context: this will not be a fairy tale



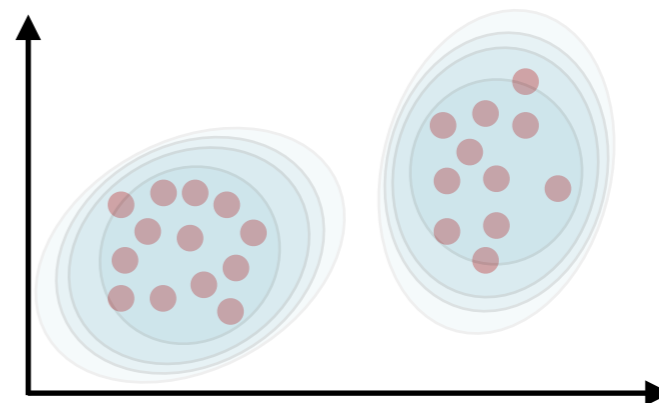
θ Parameters that solves a specific tasks

Example: K-means



Example: GMM fitting

$$\theta = \{\alpha_k, \mu_k, \Sigma_k\}_{k \in [K]}$$



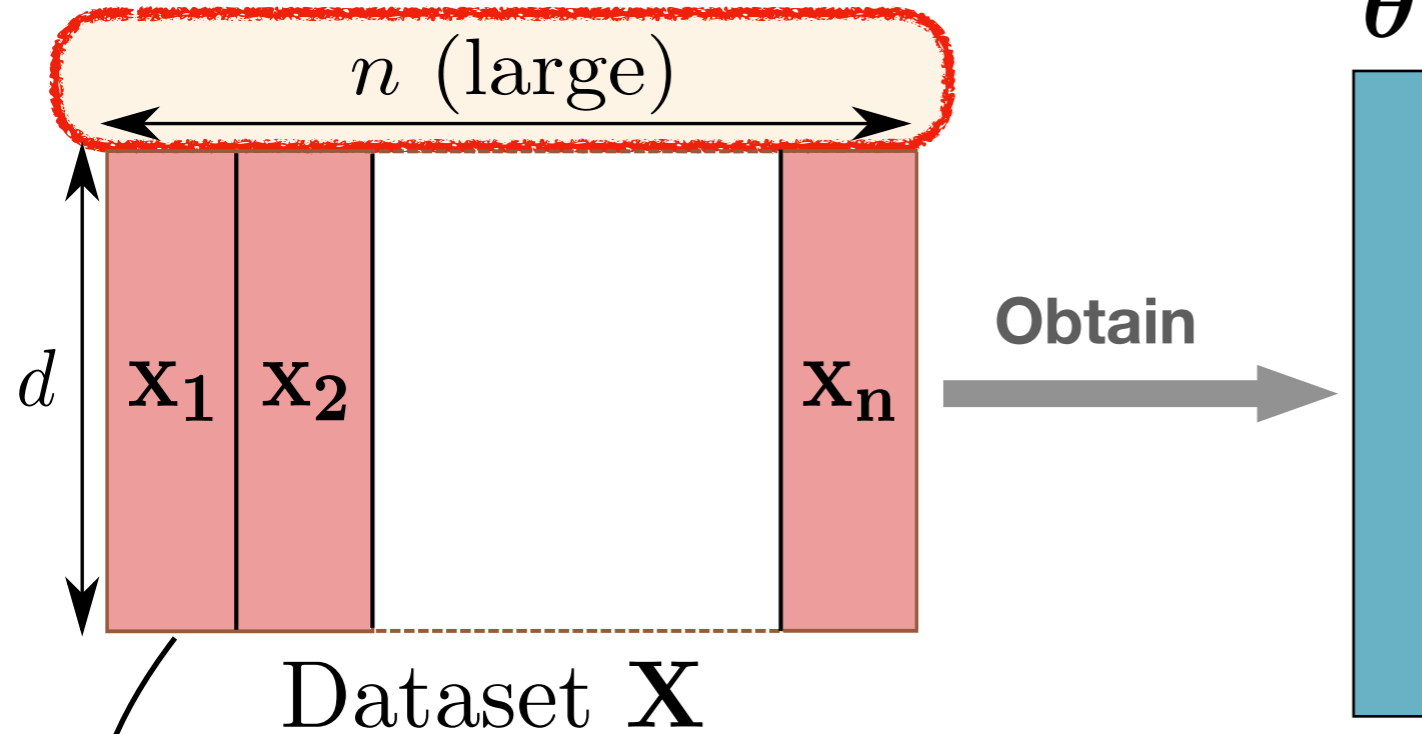
a sample (e.g. vector, image, embedding of words)



we are not like John

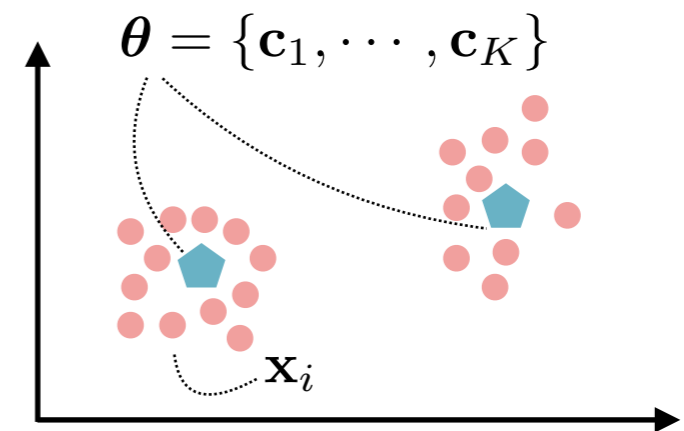
Motivations of this talk

Context: this will not be a fairy tale



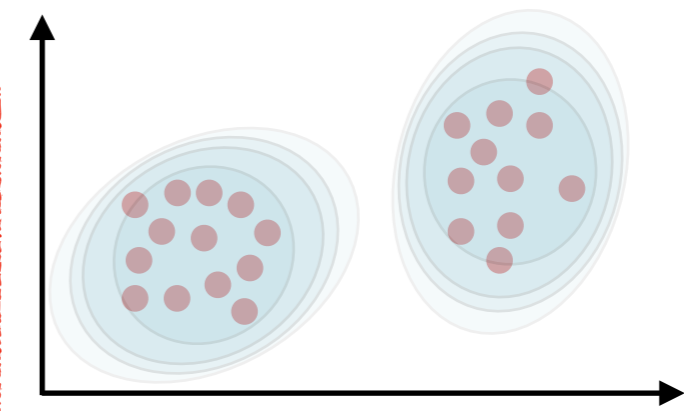
θ Parameters that solves a specific tasks

Example: K-means



Example: GMM fitting

$$\theta = \{\alpha_k, \mu_k, \Sigma_k\}_{k \in [K]}$$



a sample (e.g. vector, image, embedding of words)

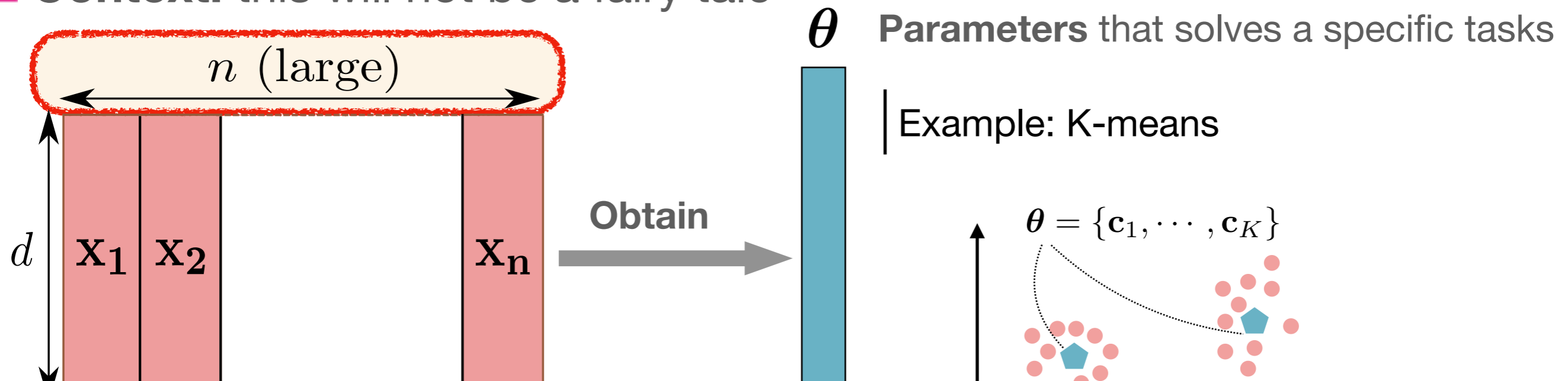
we are not like John



Large scale Machine Learning

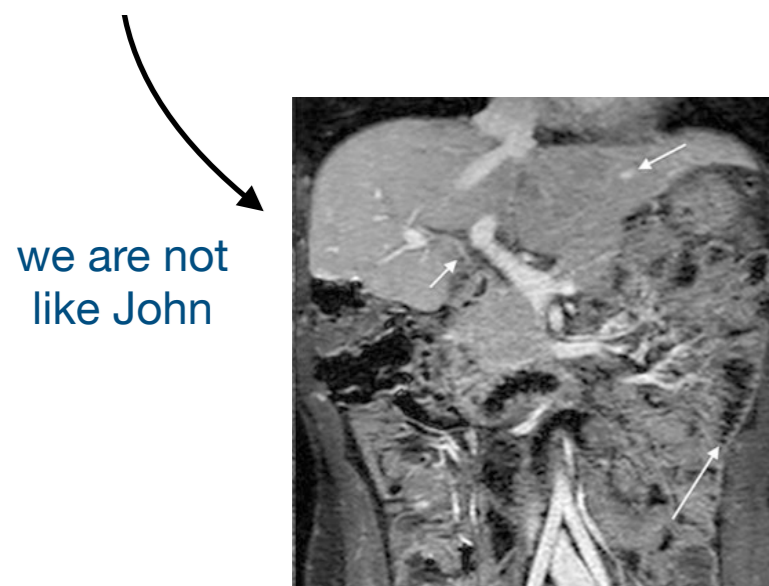
Motivations of this talk

Context: this will not be a fairy tale

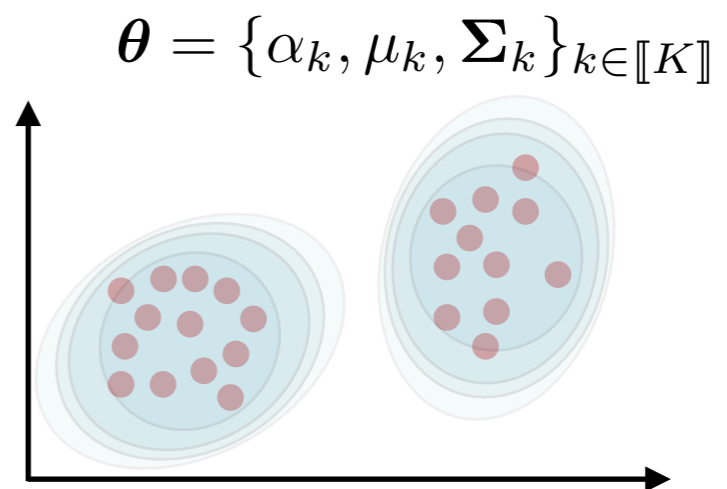


Energy and Policy Considerations for Deep Learning in NLP

Emma Strubell, Ananya Ganesh, Andrew McCallum



**Large scale
Machine Learning**



| Overview of the talk

- Part I: **A journey in the compressed sensing theory**

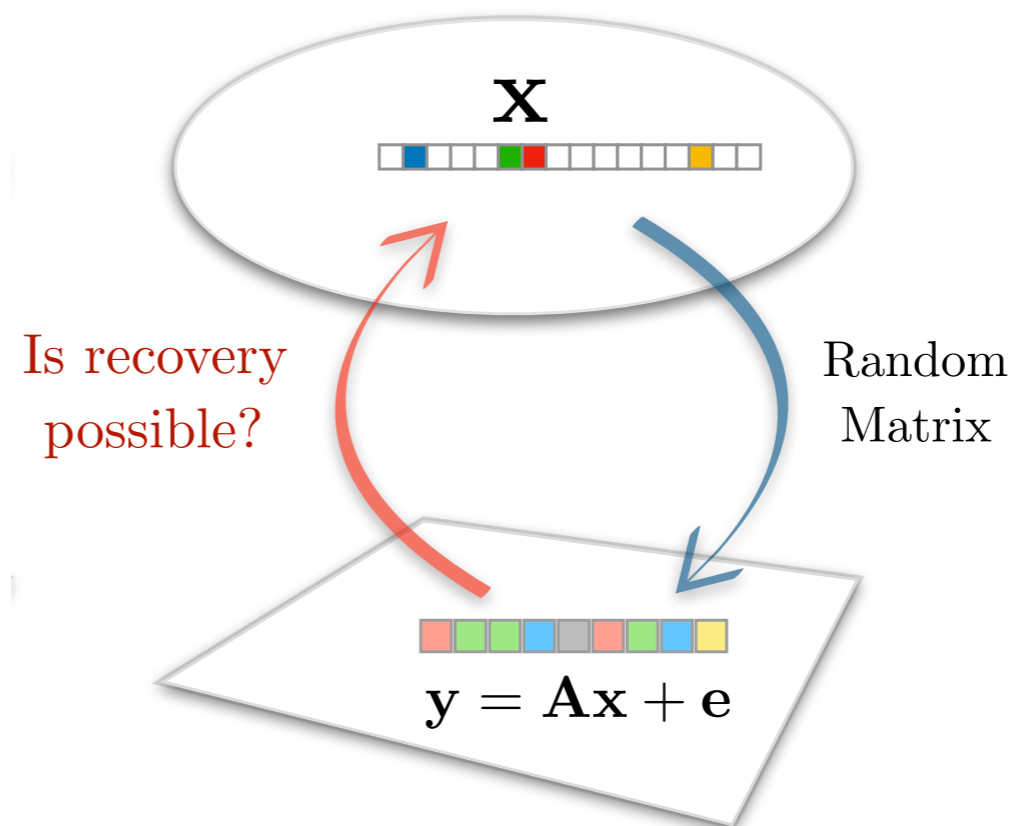
- Part II: **A bit of machine learning theory**

- Part III: **The sketching approach**

 - Applied sketching

 - Theoretical guarantees

A journey in the compressed sensing theory



| Compressed sensing theory: an invitation

■ From the basements:

- A signal: $\mathbf{x} \in \mathbb{R}^d$
- An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

Observation

$$\mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^m$$

■ Goal: recover $\mathbf{x} \in \mathbb{R}^d$ from $\mathbf{y} \in \mathbb{R}^m$

Compressed sensing theory: an invitation

From the basements:

■ A signal: $\mathbf{x} \in \mathbb{R}^d$

■ An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

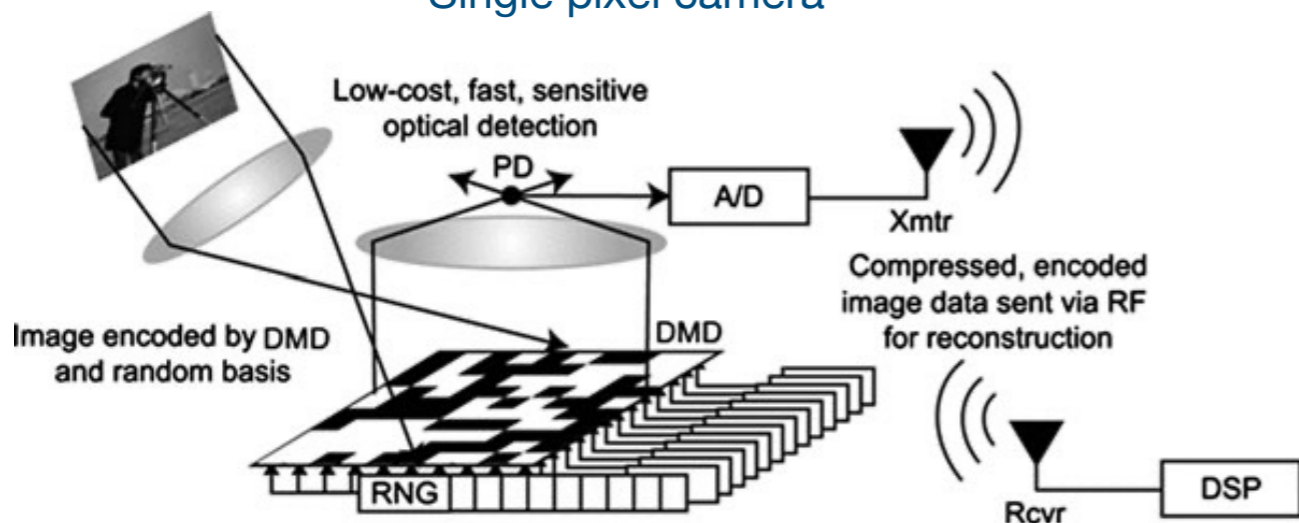
Observation

$$\mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^m$$

Goal: recover $\mathbf{x} \in \mathbb{R}^d$ from $\mathbf{y} \in \mathbb{R}^m$

■ Basis of most devices: analog-to-digital, medical imaging, radar, mobile

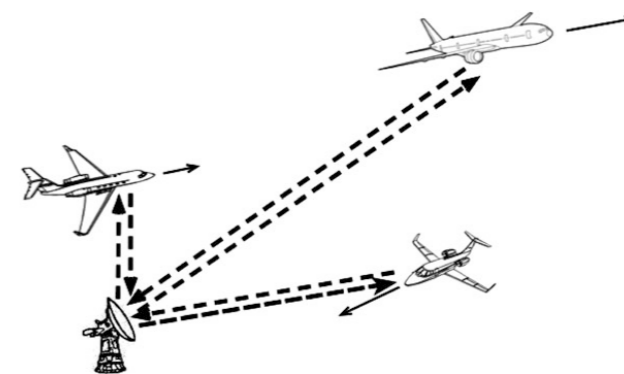
Single pixel camera



MRI



Radar



| Compressed sensing theory: an invitation

■ A signal: $\mathbf{x} \in \mathbb{R}^d$ ■ An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

Linear system: Find $\mathbf{x} \in \mathbb{R}^d$ s.t. $\mathbf{y} = \mathbf{A}\mathbf{x}$

■ What does the theory says ? $\emptyset, 1, \infty$

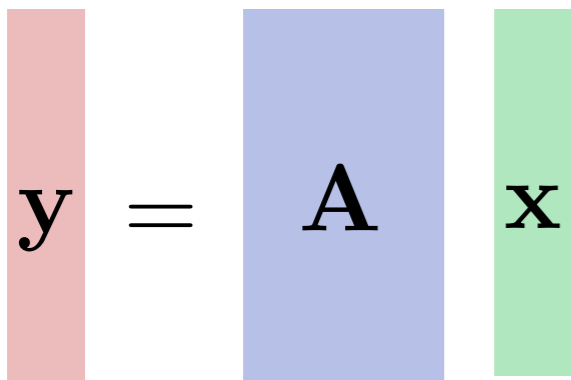
Compressed sensing theory: an invitation

■ A signal: $\mathbf{x} \in \mathbb{R}^d$ ■ An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

Linear system: Find $\mathbf{x} \in \mathbb{R}^d$ s.t. $\mathbf{y} = \mathbf{A}\mathbf{x}$

■ What does the theory says ? $\emptyset, 1, \infty$

Determined $m = d$



The diagram illustrates the equation $\mathbf{y} = \mathbf{A}\mathbf{x}$ using three vertical bars. The first bar is red and contains the variable \mathbf{y} . The second bar is blue and contains the matrix \mathbf{A} . The third bar is green and contains the variable \mathbf{x} . The bars are arranged horizontally with an equals sign between the first and second bars, and between the second and third bars.

As much data as the
size of the signal

If well behaved **unique sol**

$$\mathbf{x} := \mathbf{A}^{-1}\mathbf{y}$$

Compressed sensing theory: an invitation

■ A signal: $\mathbf{x} \in \mathbb{R}^d$ ■ An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

Linear system: Find $\mathbf{x} \in \mathbb{R}^d$ s.t. $\mathbf{y} = \mathbf{A}\mathbf{x}$

■ What does the theory says ? $\emptyset, 1, \infty$

Determined $m = d$

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

| As much data as the size of the signal

| If well behaved **unique sol**

$$\mathbf{x} := \mathbf{A}^{-1}\mathbf{y}$$

Overdetermined $m > d$

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

| We have more data

| From **low dim** to **high dim**

| In general **no sol**:

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

[Gauss, 1795]

Compressed sensing theory: an invitation

■ A signal: $\mathbf{x} \in \mathbb{R}^d$ ■ An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

Linear system: Find $\mathbf{x} \in \mathbb{R}^d$ s.t. $\mathbf{y} = \mathbf{A}\mathbf{x}$

■ What does the theory says ? $\emptyset, 1, \infty$

Determined $m = d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

As much data as the size of the signal

If well behaved **unique sol**

$$\mathbf{x} := \mathbf{A}^{-1} \mathbf{y}$$

Overdetermined $m > d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

We have more data

From **low dim** to **high dim**

In general **no sol**:

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

[Gauss, 1795]

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

We did **compression**

From **high dim** to **low dim**

Infinity of solutions

$$\mathbf{x}_0 + \lambda \mathbf{z}, \mathbf{z} \in \ker(\mathbf{A})$$

Compressed sensing theory: an invitation

■ A signal: $\mathbf{x} \in \mathbb{R}^d$ ■ An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

Linear system: Find $\mathbf{x} \in \mathbb{R}^d$ s.t. $\mathbf{y} = \mathbf{A}\mathbf{x}$

■ What does the theory says ? $\emptyset, 1, \infty$

Determined $m = d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

As much data as the size of the signal

If well behaved **unique sol**

$$\mathbf{x} := \mathbf{A}^{-1} \mathbf{y}$$

Overdetermined $m > d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

We have more data

From **low dim** to **high dim**

In general **no sol**:

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

[Gauss, 1795]

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

We did **compression**

From **high dim** to **low dim**

Infinity of solutions

$$\mathbf{x}_0 + \lambda \mathbf{z}, \mathbf{z} \in \ker(\mathbf{A})$$

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

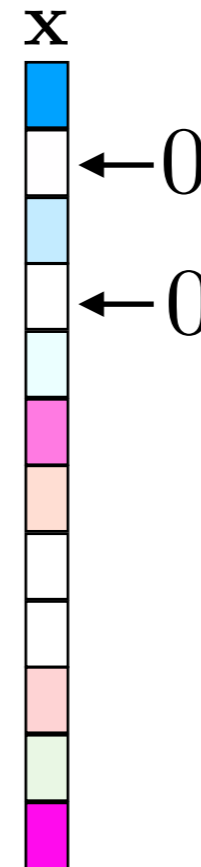
Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The sparse assumption

Recovery is possible when we **know something more** about \mathbf{x}

The signal $\mathbf{x} \in \mathbb{R}^d$ is high-dim **but it is full of zeros**



Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The sparse assumption

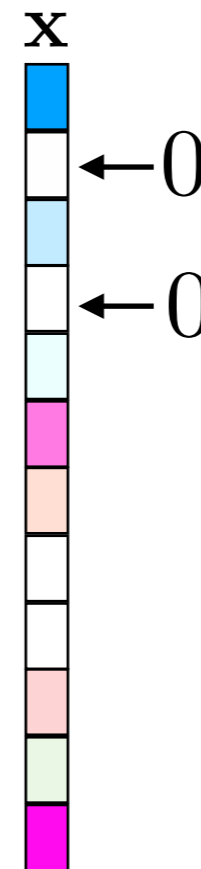
Recovery is possible when we **know something more** about \mathbf{x}

The signal $\mathbf{x} \in \mathbb{R}^d$ is high-dim **but it is full of zeros**

In a way \mathbf{x} **lives in a low-dim. space**

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

We need much less information to recover \mathbf{x}



Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The sparse assumption

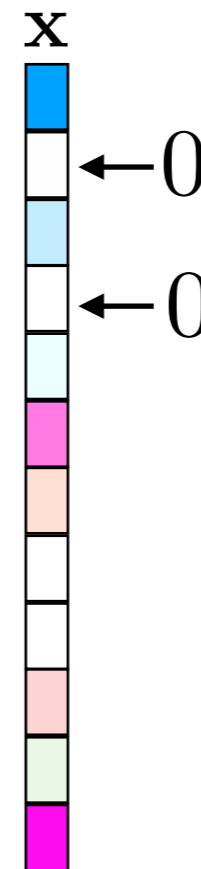
Recovery is possible when we **know something more** about \mathbf{x}

The signal $\mathbf{x} \in \mathbb{R}^d$ is high-dim **but it is full of zeros**

In a way \mathbf{x} **lives in a low-dim. space**

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

We need much less information to recover \mathbf{x}



Number of nnz components

$$\|\mathbf{x}\|_0 = \sum_{i=1}^d \mathbf{1}_{x_i \neq 0}$$

E.g. $\|\mathbf{x}\|_0 \leq k$

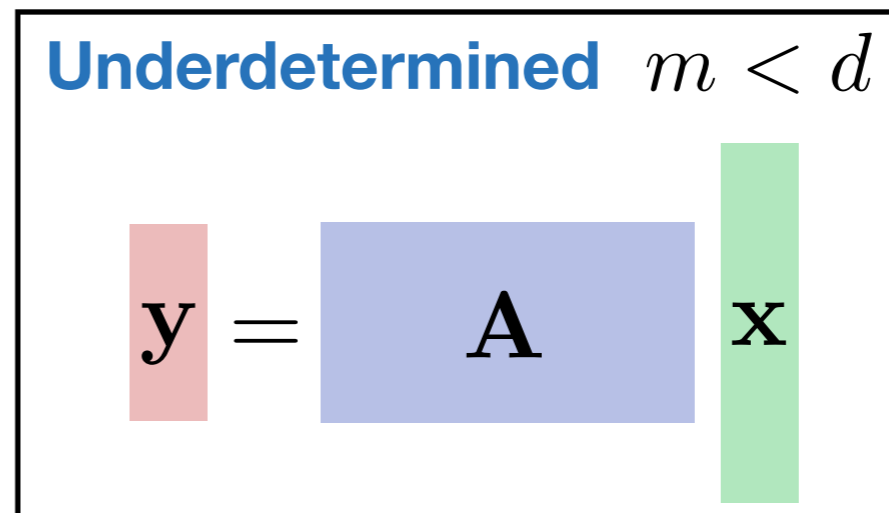
\mathbf{x} is k -sparse

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions



The sparse assumption

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

We could solve:

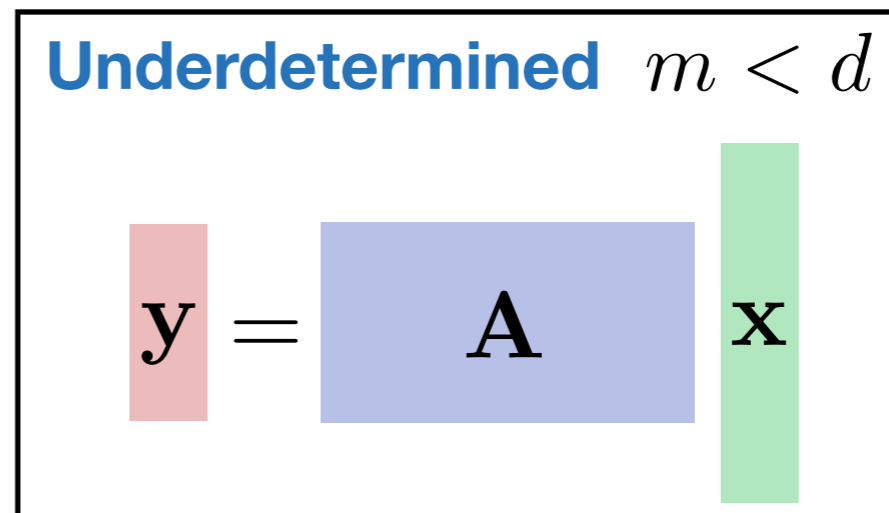
$$(1) \quad \min_{\mathbf{x} \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}} \|\mathbf{x}\|_0$$

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions



The sparse assumption

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

We could solve:

$$(1) \quad \min_{\mathbf{x} \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}} \|\mathbf{x}\|_0 \longrightarrow \text{Find the sparsest vector}$$

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The sparse assumption

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

We could solve:

$$(1) \quad \min_{\mathbf{x} \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}} \|\mathbf{x}\|_0$$

Which satisfies the constraints

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The sparse assumption

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

We could solve:

$$(1) \quad \min_{\mathbf{x} \text{ s.t. } \mathbf{y} = \mathbf{A} \mathbf{x}} \|\mathbf{x}\|_0$$

Theorem $d \geq 2k$

$\exists \mathbf{A} \in \mathbb{R}^{2k \times d}$ ($m = 2k$)

every k -sparse \mathbf{x}

can be recovered from $\mathbf{y} = \mathbf{A} \mathbf{x}$

by solving (1)

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The sparse assumption

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

We could solve:

$$(1) \quad \min_{\mathbf{x} \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}} \|\mathbf{x}\|_0$$

Theorem $d \geq 2k$

$\exists \mathbf{A} \in \mathbb{R}^{2k \times d}$ ($m = 2k$)

every k -sparse \mathbf{x}

can be recovered from $\mathbf{y} = \mathbf{A}\mathbf{x}$

by solving (1)

Is it done ? **No ...**

Not robust to noise

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

(1) is NP-hard

Not robust w.r.t. sparsity level

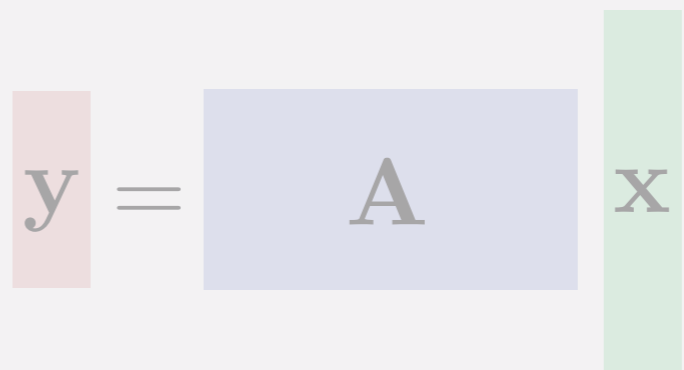
Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$



The sparse solution

There is a unique

with the minimum

$$(1) \min_{\mathbf{x} \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}} \|\mathbf{x}\|_0$$

Is \mathbf{x} sparse in practice ??

by solving (1)

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

Is it done ? **No ...**

Not robust to noise

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

(1) is NP-hard

Not robust w.r.t. sparsity level

| Compressed sensing theory: an invitation

- Most data encountered are sparse in **another representation**

Original signal $\longrightarrow \mathbf{x} = \mathbf{D}\mathbf{x}_0 \longleftarrow$ sparse vector

- $\mathbf{D} \in \mathbb{R}^{d \times d}$ is another **ortho. basis** than the canonical one $\mathbf{D}^{-1} = \mathbf{D}^\top$

$$\mathbf{x} = \mathbf{D}\mathbf{x}_0 \iff \mathbf{x}_0 = \mathbf{D}^\top \mathbf{x}$$

Compressed sensing theory: an invitation

Most data encountered are sparse in **another representation**

Original signal $\longrightarrow \mathbf{x} = \mathbf{D}\mathbf{x}_0 \longleftarrow$ sparse vector

$\mathbf{D} \in \mathbb{R}^{d \times d}$ is another **ortho. basis** than the canonical one $\mathbf{D}^{-1} = \mathbf{D}^\top$

Discrete Fourier basis (FFT)

Discrete Cosinus Transform

Wavelet

$$D_{pq} = \exp\left(-2i\frac{\pi}{d}pq\right)$$

$$D_{pq} = \cos\left(\frac{\pi}{d}\left(p + \frac{1}{2}\right)q\right)$$

[Haar, 1909]

[Gabor, 1946]

[Morlet & Grossmann, 1984]

[Mallat, 1986]

[Daubechies, 1987]



Compressed sensing theory: an invitation

Most data encountered are sparse in **another representation**

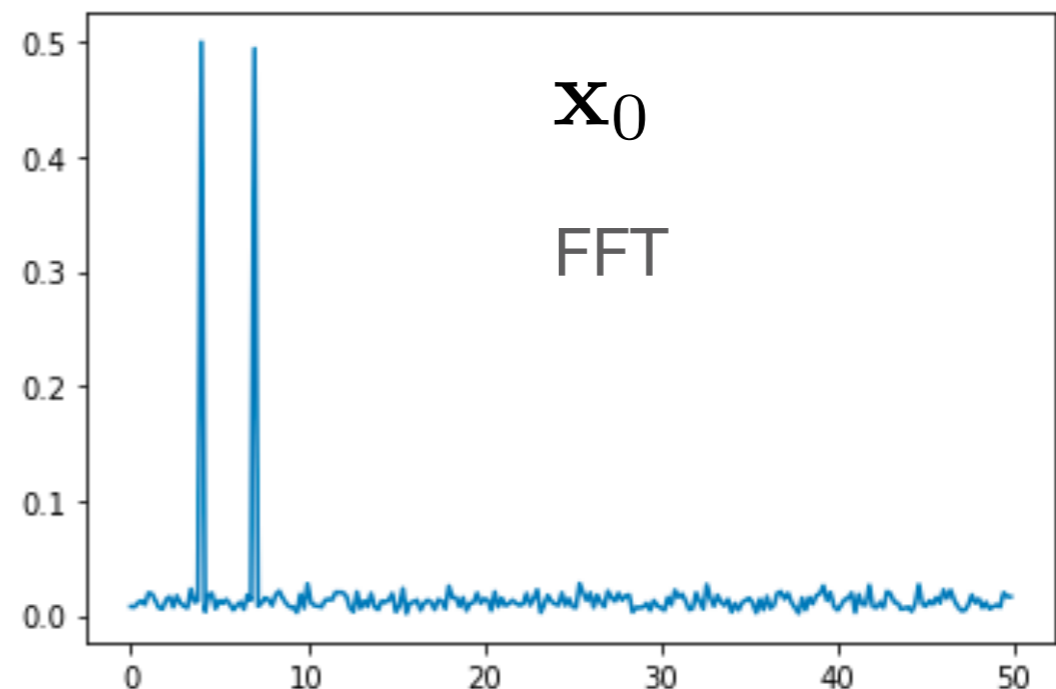
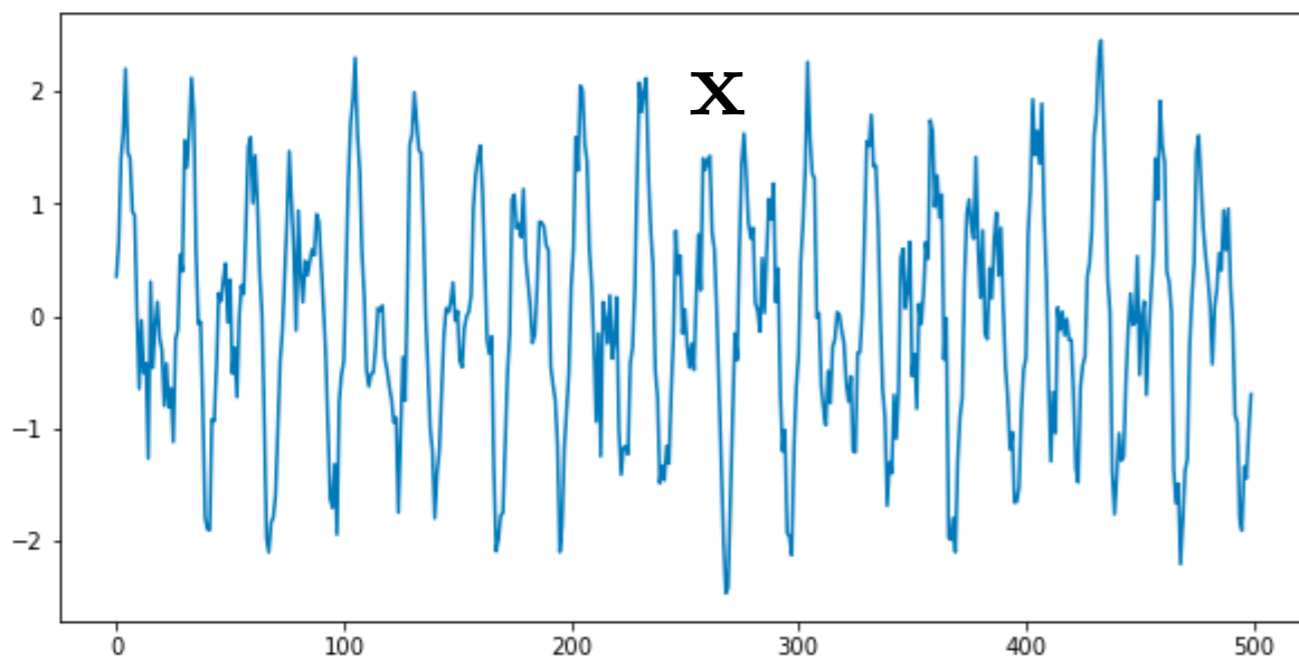
Original signal $\longrightarrow \mathbf{x} = \mathbf{D}\mathbf{x}_0 \longleftarrow$ sparse vector

$\mathbf{D} \in \mathbb{R}^{d \times d}$ is another **ortho. basis** than the canonical one $\mathbf{D}^{-1} = \mathbf{D}^\top$

Discrete Fourier basis (FFT) Discrete Cosinus Transform Wavelet

$$D_{pq} = \exp(-2i\frac{\pi}{d}pq) \quad D_{pq} = \cos(\frac{\pi}{d}(p + \frac{1}{2})q)$$

Quite magical ...



Compressed sensing theory: an invitation

Most data encountered are sparse in **another representation**

Original signal $\longrightarrow \mathbf{x} = \mathbf{D}\mathbf{x}_0 \longleftarrow$ sparse vector

$\mathbf{D} \in \mathbb{R}^{d \times d}$ is another **ortho. basis** than the canonical one $\mathbf{D}^{-1} = \mathbf{D}^\top$

Discrete Fourier basis (FFT)

Discrete Cosinus Transform

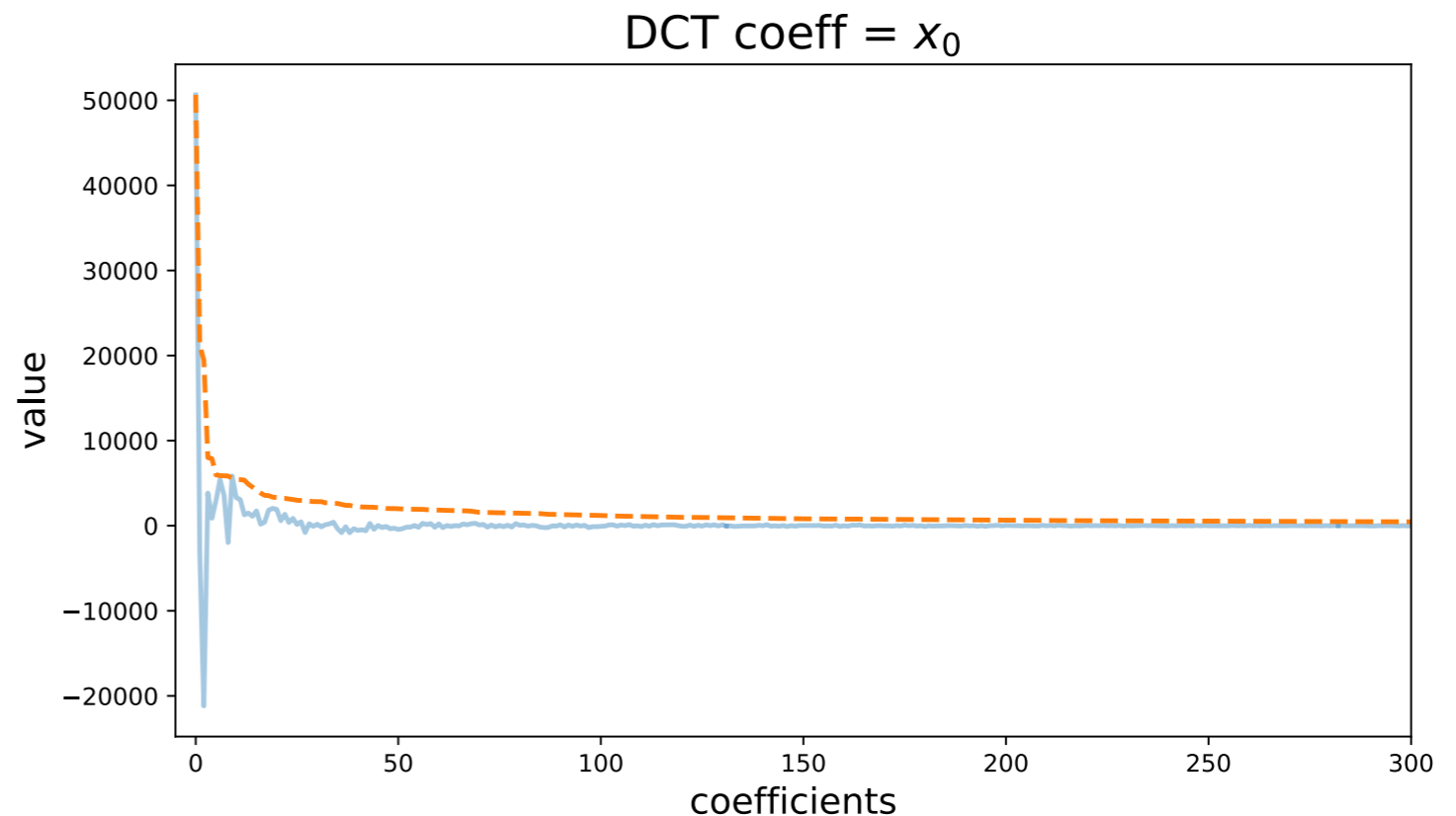
Wavelet

$$D_{pq} = \exp\left(-2i\frac{\pi}{d}pq\right)$$

$$D_{pq} = \cos\left(\frac{\pi}{d}\left(p + \frac{1}{2}\right)q\right)$$

Quite magical ...

$$d = 512 * 512$$



| Compressed sensing theory: an invitation

Most data encountered are sparse in **another representation**

Original signal $\longrightarrow \mathbf{x} = \mathbf{D}\mathbf{x}_0 \longleftarrow$ sparse vector

$\mathbf{D} \in \mathbb{R}^{d \times d}$ is another **ortho. basis** than the canonical one $\mathbf{D}^{-1} = \mathbf{D}^\top$

Discrete Fourier basis (FFT)

Discrete Cosinus Transform

Wavelet

$$D_{pq} = \exp\left(-2i\frac{\pi}{d}pq\right)$$

$$D_{pq} = \cos\left(\frac{\pi}{d}\left(p + \frac{1}{2}\right)q\right)$$

At the core of many **compression schemes** (JPEG, MPEG, MP3 ...)

Keep only largest value of \mathbf{x}_0

Wavelet compression

Original

Ratio = 1.5%
PSNR= 29.31dB

Ratio = 3%
PSNR= 33.62dB

Ratio = 5%
PSNR= 37.81dB

Ratio = 10%
PSNR= 44.62dB

Ratio = 50%
PSNR= 66.89dB



Compressed sensing theory: an invitation

Most data encountered are sparse in **another representation**

Original signal $\longrightarrow \mathbf{x} = \mathbf{D}\mathbf{x}_0 \longleftarrow$ sparse vector

$\mathbf{D} \in \mathbb{R}^{n \times n}$ is another **ortho. basis** than the canonical one $\mathbf{D}^{-1} = \mathbf{D}^\top$

Discret

$D_{pq} = e$

\mathbf{x} is reasonably sparse if we look right

At the c

Wavelet

..)

Keep only largest value of \mathbf{x}_0

Wavelet compression

Original

Ratio = 1.5%
PSNR= 29.31dB

Ratio = 3%
PSNR= 33.62dB

Ratio = 5%
PSNR= 37.81dB

Ratio = 10%
PSNR= 44.62dB

Ratio = 50%
PSNR= 66.89dB

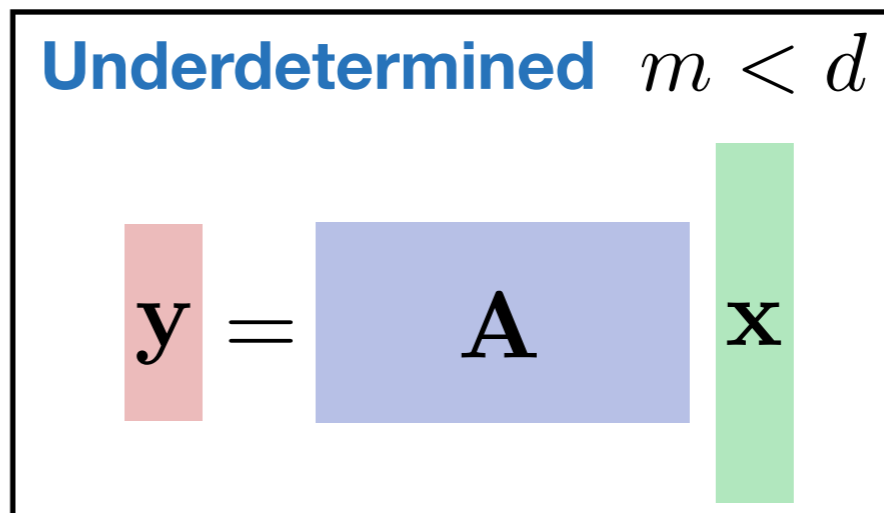


Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions



The sparse assumption

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**



We could solve: (1) $\min_{\mathbf{x} \text{ s.t. } \mathbf{y}=\mathbf{A}\mathbf{x}} \|\mathbf{x}\|_0$

Is it done ? **No ...**

Not robust to noise

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

(1) is NP-hard

Not robust w.r.t. sparsity level

Compressed sensing theory: an invitation

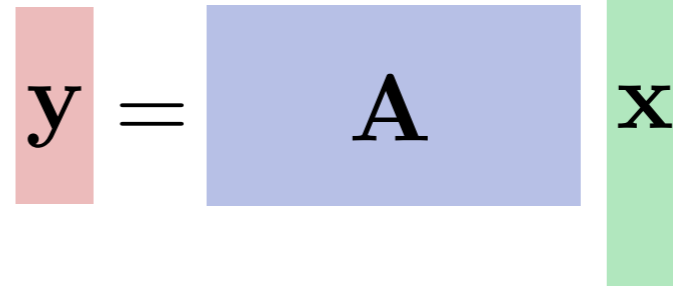
How can we recover \mathbf{x} ?

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

The **LASSO**

[Thibshirani, 1996]
[Chen & Donoho, 1995]

Underdetermined $m < d$


$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

$$\min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{A} \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

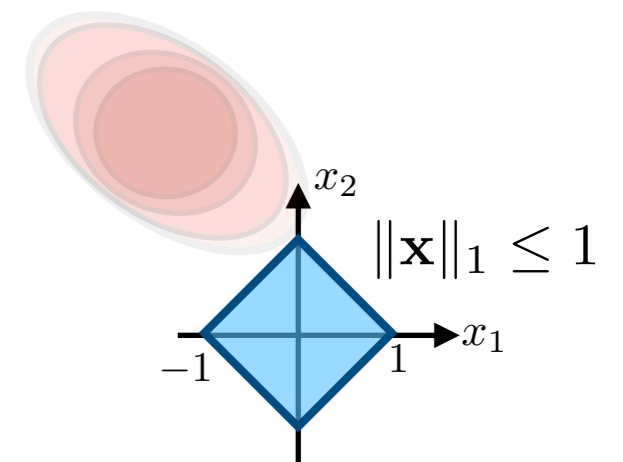
The LASSO

[Thibshirani, 1996]
[Chen & Donoho, 1995]

$$\min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

$\mathbf{A}\mathbf{x}$ is closed to \mathbf{y}

\mathbf{x} is penalized so as to have small L1 norm
promotes sparsity



Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

- The signal $\mathbf{x} \in \mathbb{R}^d$ is **sparse**

Underdetermined $m < d$

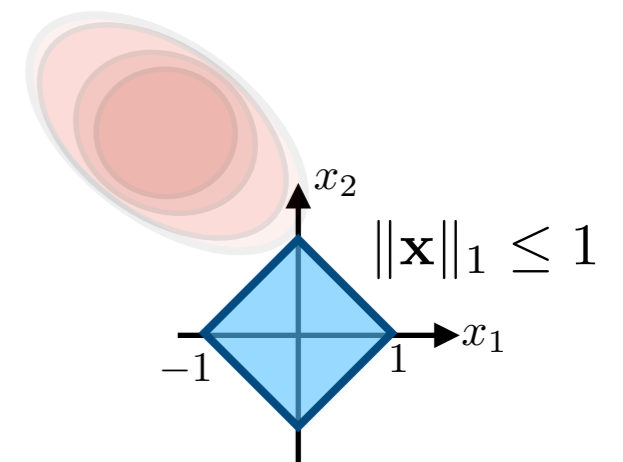
$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The LASSO

- Major impact on ML
[Bach & al, 2012]

$$\min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

- Strictly convex



Can be solved using many algorithms...

- Proximal algorithms (iterative thresholding): ISTA, FISTA

[Beck & Teboulle, 2009]

- (Block) Coordinate descent algorithms [Friedman & al, 2007]

- Here we will use CELER [Massias & al, 2018]



| Compressed sensing theory: an invitation

■ The LASSO in practice for CS



| Compressed sensing theory: an invitation

■ The LASSO in practice for CS

$$d = 128 * 128$$

X



Compressed sensing theory: an invitation

The LASSO in practice for CS

$$d = 128 * 128$$

X



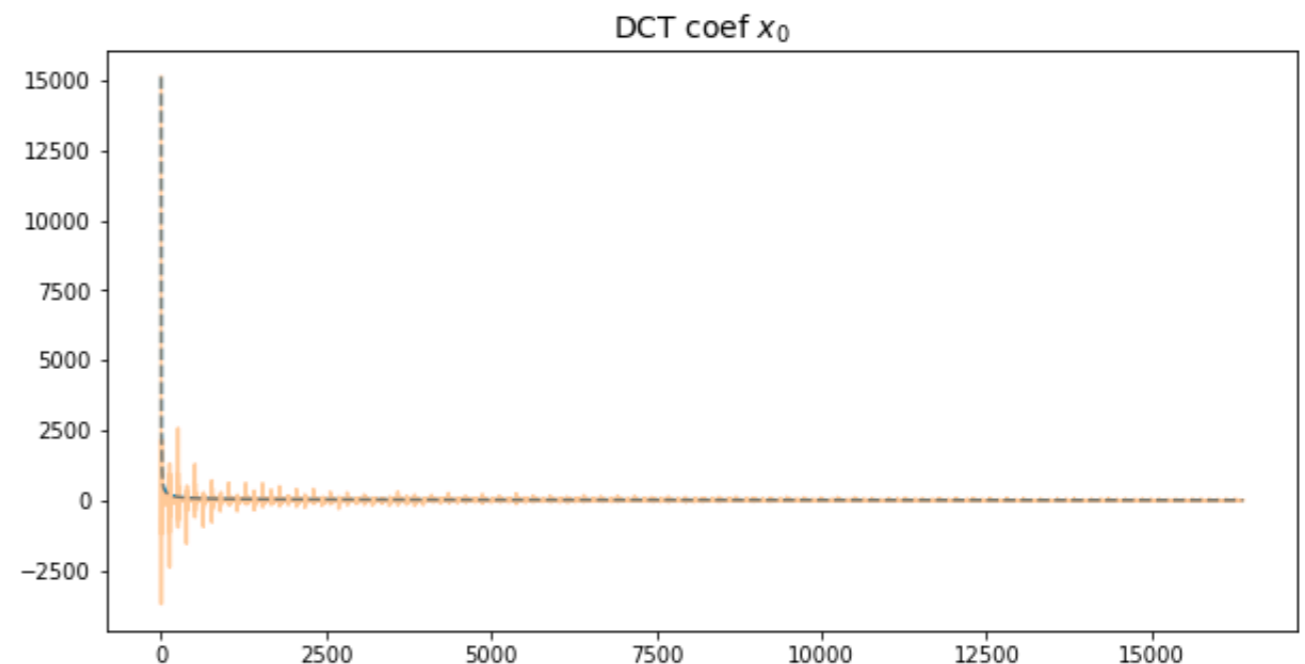
■ An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

| We do **compression** $m < d$

| We take $A_{ij} \sim \frac{1}{\sqrt{m}} \mathcal{N}(0, 1)$

Observation $\mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^m$

■ **X** is sparse in another basis (DCT)



Compressed sensing theory: an invitation

The LASSO in practice for CS

$$d = 128 * 128$$

X



- An acquisition system: $\mathbf{A} \in \mathbb{R}^{m \times d}$

| We do **compression** $m < d$

| We take $A_{ij} \sim \frac{1}{\sqrt{m}} \mathcal{N}(0, 1)$

Observation $\mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^m$

- **X** is sparse in another basis (DCT)
- We solve the LASSO

$$\min_{\mathbf{x}_0 \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{A}\mathbf{D}\mathbf{x}_0\|_2^2 + \lambda \|\mathbf{x}_0\|_1$$

Compressed sensing theory: an invitation

The LASSO in practice for CS

$$m = 70\%d$$

Original



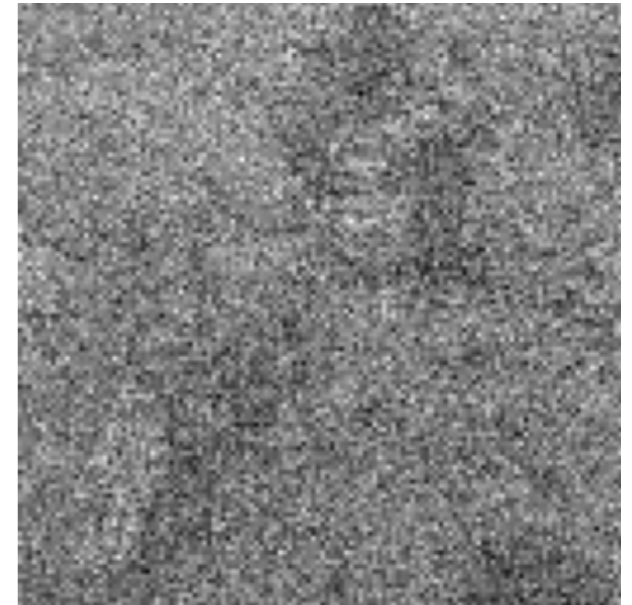
LASSO



OMP



Linear reg.



Original



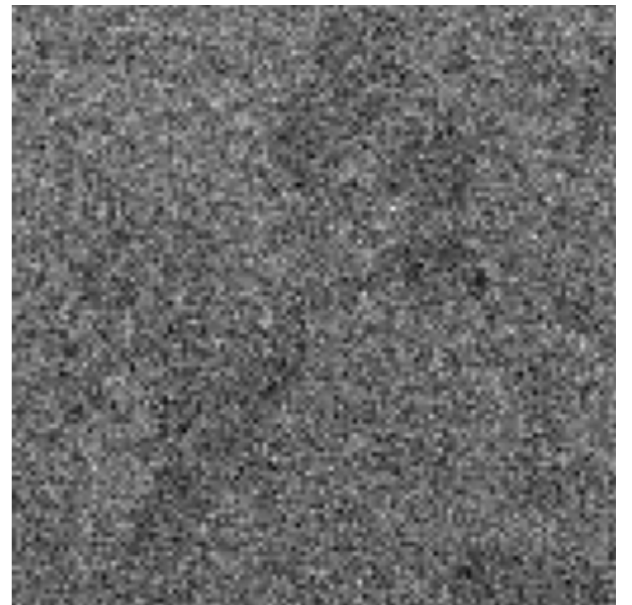
LASSO



OMP

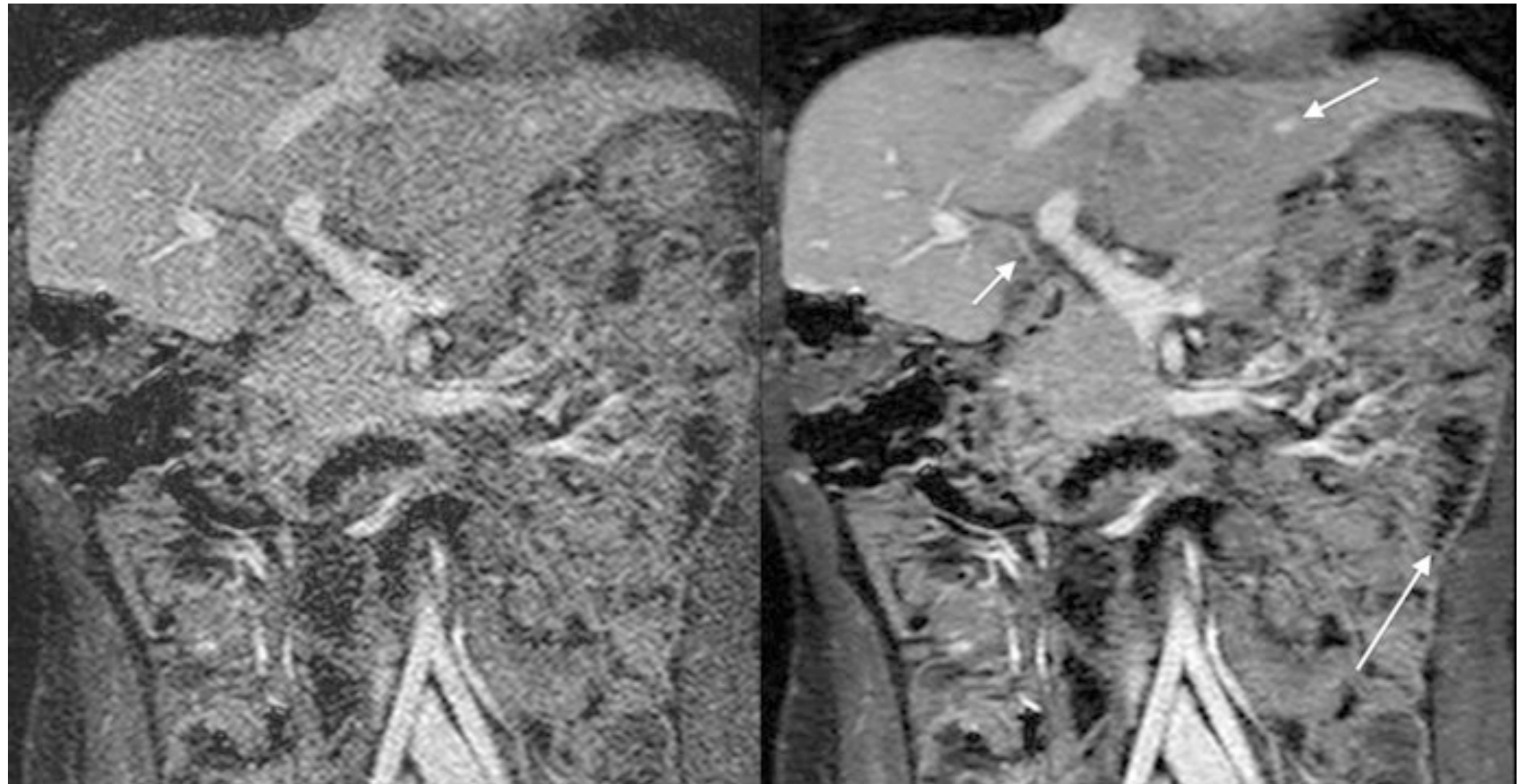


Linear reg.



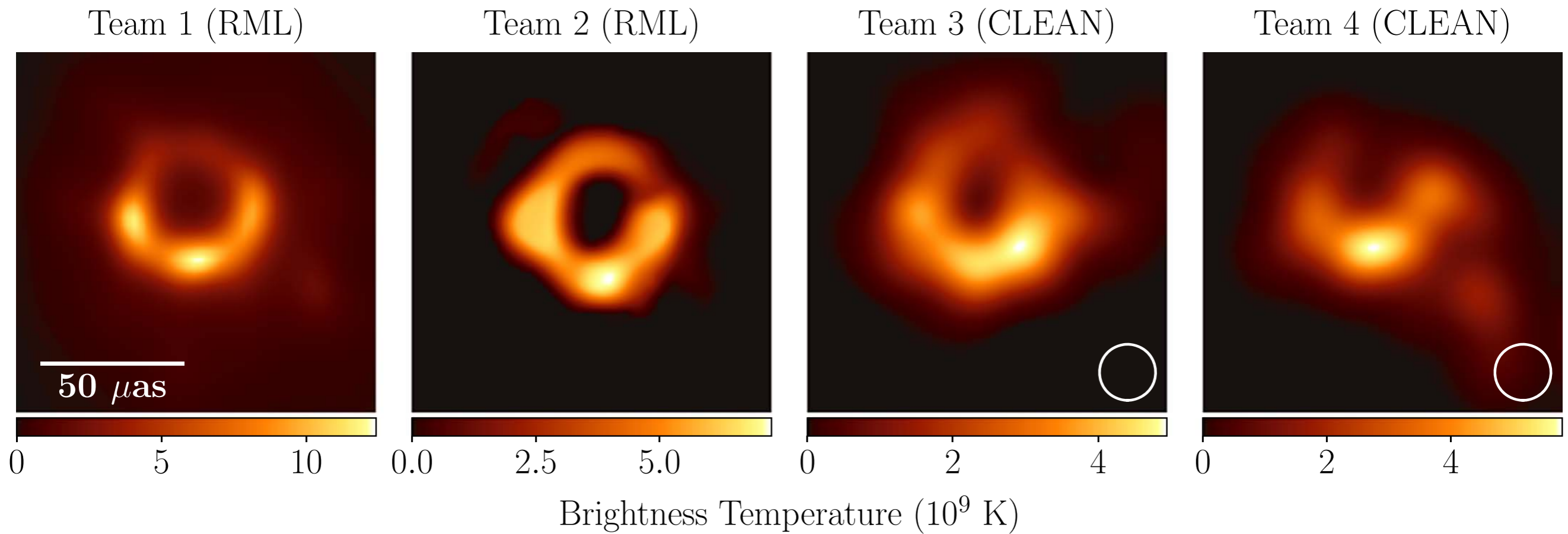
| Compressed sensing theory: an invitation

■ Other CS results



Compressed sensing theory: an invitation

Other CS results



| Compressed sensing theory: an invitation

■ Guarantees for the LASSO

■ Does the LASSO truly recovers $\mathbf{x} \in \mathbb{R}^d$?

Compressed sensing theory: an invitation

Guarantees for the LASSO

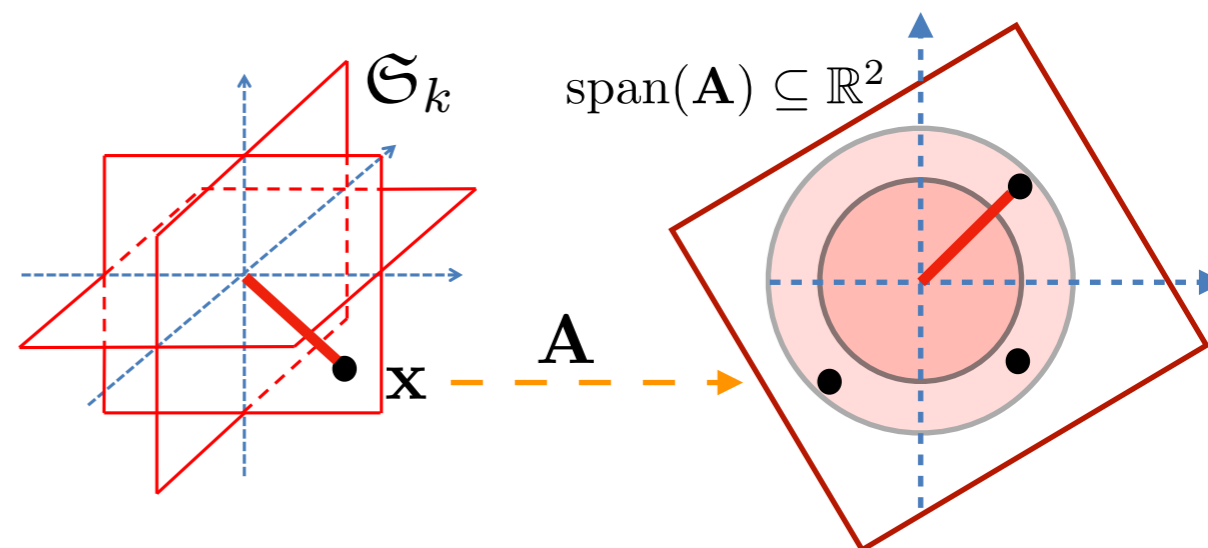
Does the LASSO truly recovers $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



Compressed sensing theory: an invitation

Guarantees for the LASSO

- Does the LASSO truly recovers $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$

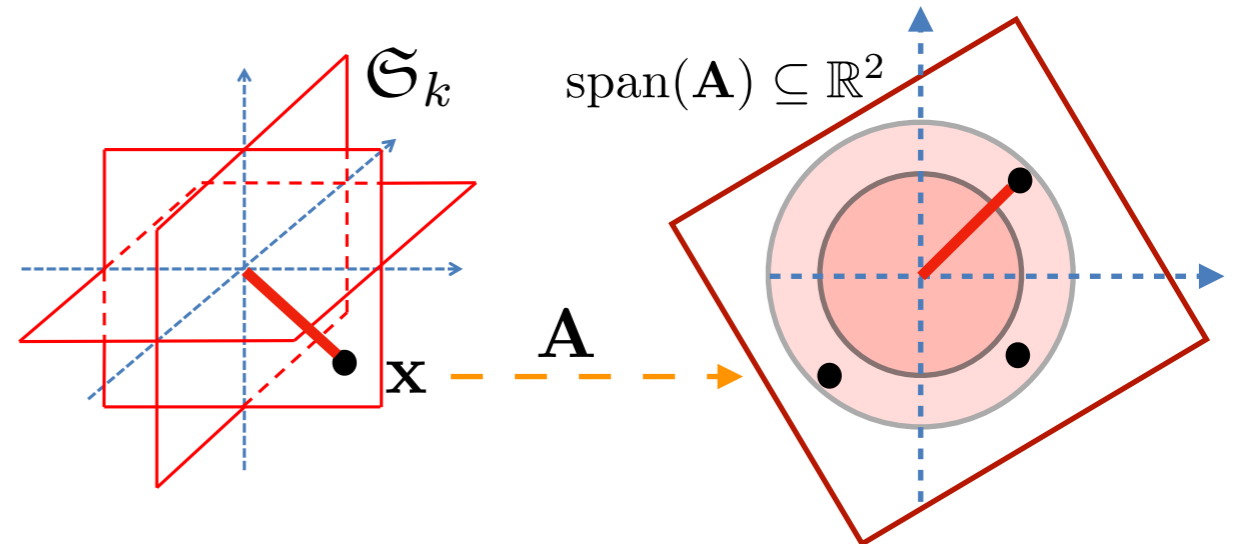
$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$

- E.g. **Gaussian matrices**

$$\mathbf{A} \in \mathbb{R}^{m \times d} \text{ with } A_{ij} \sim \frac{1}{\sqrt{m}} \mathcal{N}(0, 1)$$

$$m \gtrsim \delta^{-2} k \ln\left(e \frac{d}{k}\right)$$

with overwhelming prob. we have the RIP with δ



Compressed sensing theory: an invitation

Guarantees for the LASSO

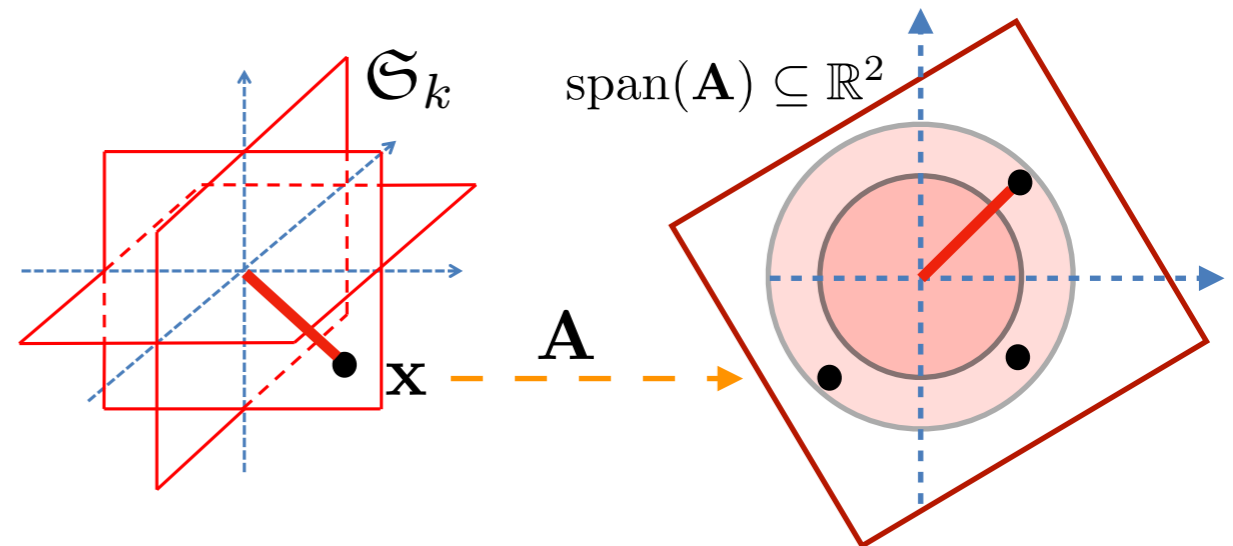
- Does the LASSO truly recovers $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



- E.g. **Gaussian matrices**

$$\mathbf{A} \in \mathbb{R}^{m \times d} \text{ with } A_{ij} \sim \frac{1}{\sqrt{m}} \mathcal{N}(0, 1)$$

$$m \gtrsim \delta^{-2} k \ln\left(e \frac{d}{k}\right)$$

with overwhelming prob. we have the **RIP** with δ

- Same holds for **AD** where **A** is **Gaussian** and **D** **orthogonal**

Compressed sensing theory: an invitation

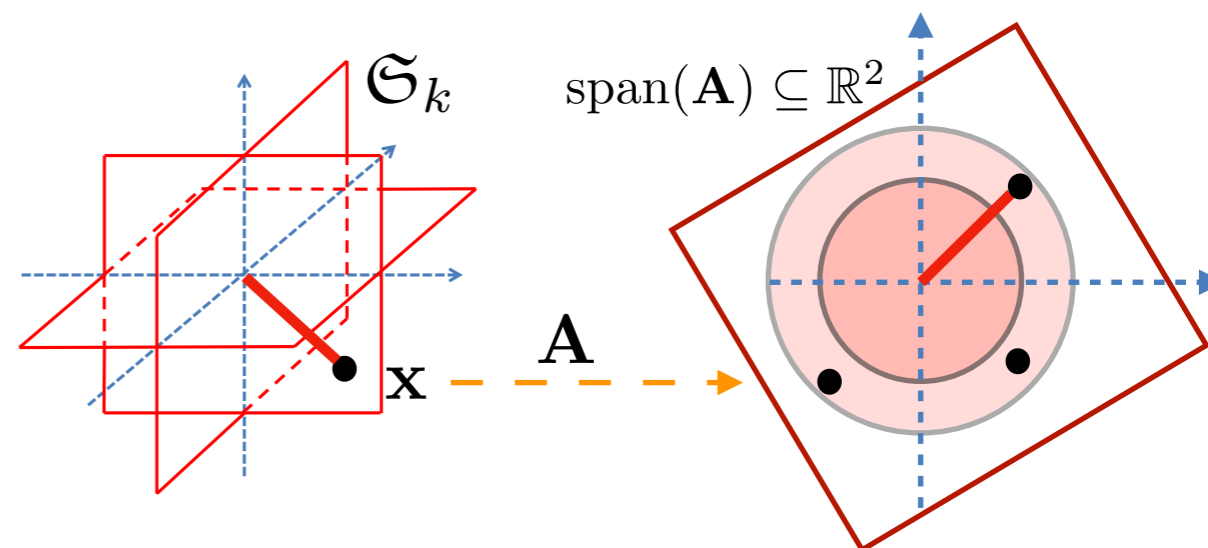
Guarantees for the LASSO

- Does the LASSO truly recovers $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$$
$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



- E.g. **Gaussian matrices**

$$m \gtrsim \delta^{-2} k \ln\left(e \frac{d}{k}\right)$$

- Related with the **magical Johnson-Lindenstrauss Lemma**

Every $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^d$ can be **linearly embedded** in \mathbb{R}^m
with $\delta > 0$ **distorsion** provided $m \gtrsim \delta^{-2} \ln(k)$

- Does not depend on d !**

Compressed sensing theory: an invitation

Guarantees for the LASSO

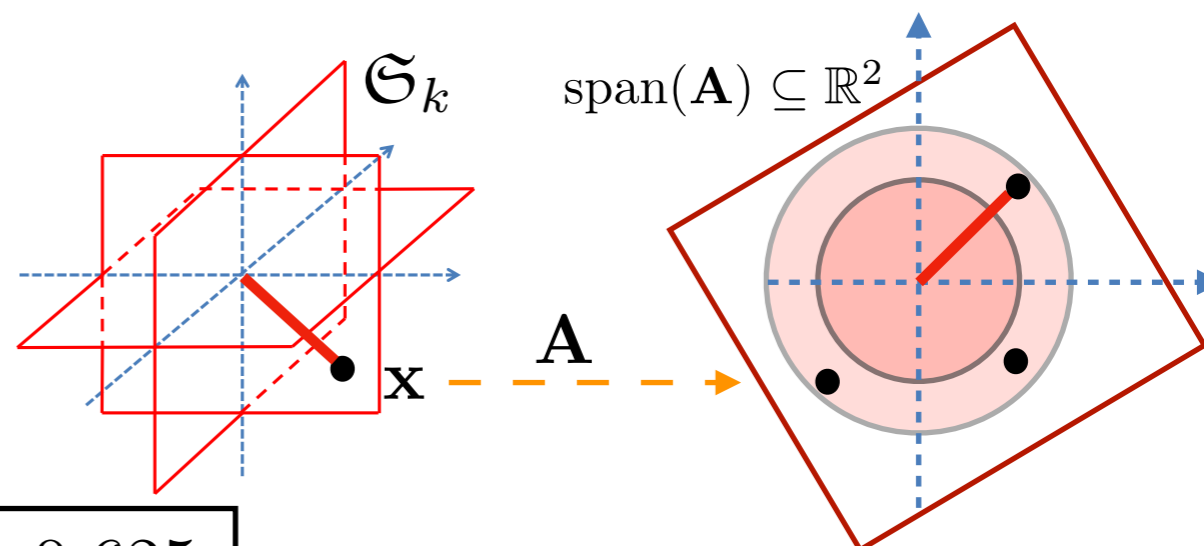
Does the LASSO truly recover $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



The result: Suppose $\delta_{2k} < 4/\sqrt{41} \approx 0.625$

Compressed sensing theory: an invitation

Guarantees for the LASSO

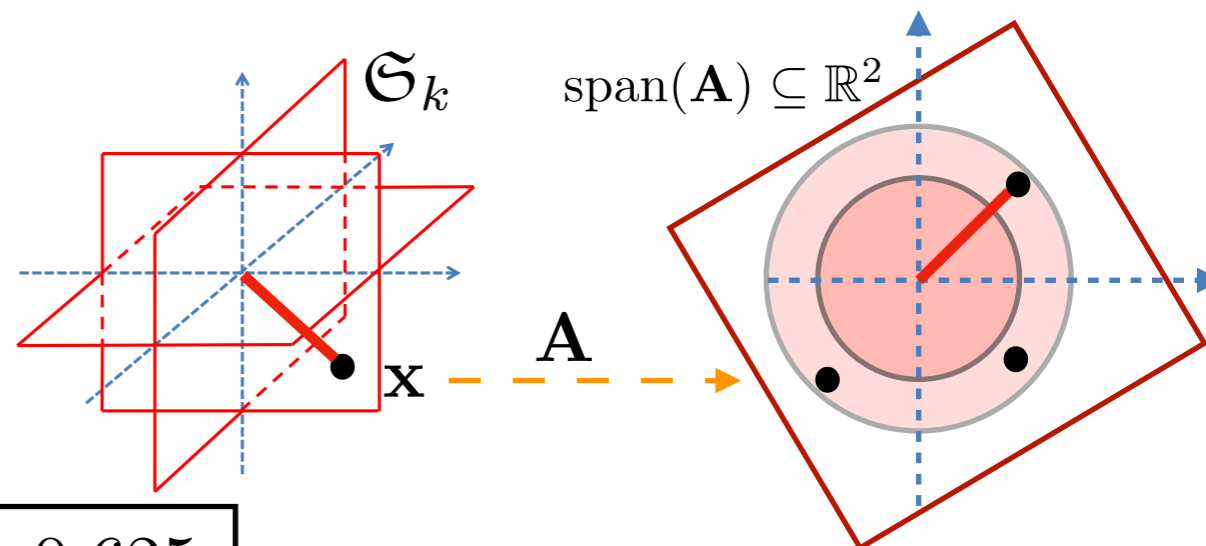
- Does the LASSO truly recover $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



The result: Suppose $\delta_{2k} < 4/\sqrt{41} \approx 0.625$

Take $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}$
with $\|\mathbf{e}\|_2 \leq \eta$

Compressed sensing theory: an invitation

Guarantees for the LASSO

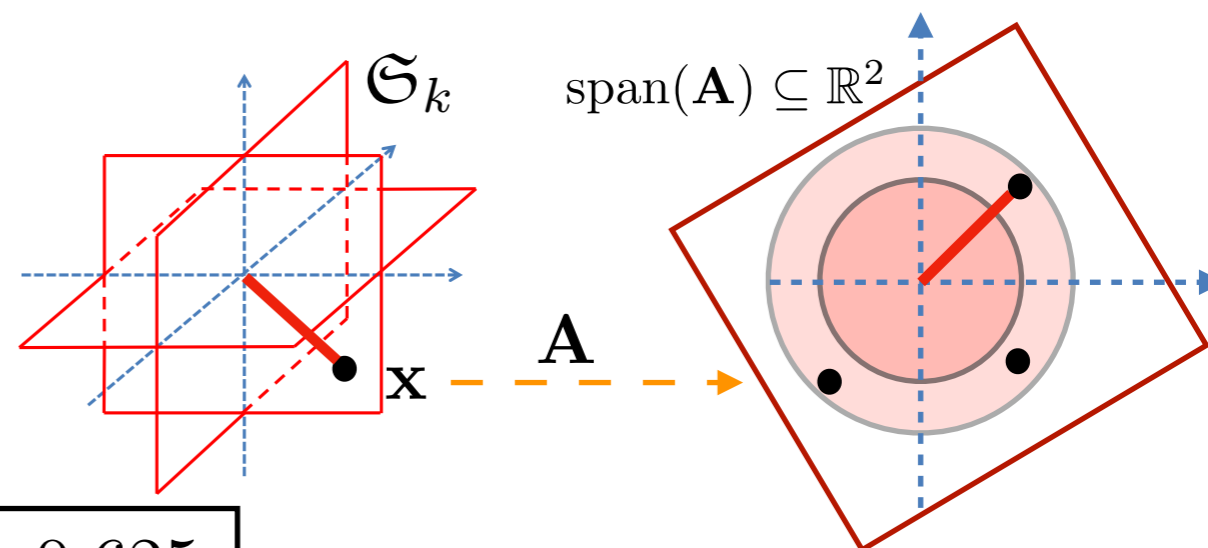
Does the LASSO truly recover $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



The result: Suppose $\delta_{2k} < 4/\sqrt{41} \approx 0.625$

Take $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} = \mathbf{Ax} + \mathbf{e}$
with $\|\mathbf{e}\|_2 \leq \eta$

$$\hat{\mathbf{x}} \text{ sol. of } \begin{cases} \min \|\mathbf{z}\|_1 \\ \text{s.t. } \|\mathbf{y} - \mathbf{Az}\|_2 \leq \eta \end{cases}$$

Compressed sensing theory: an invitation

Guarantees for the LASSO

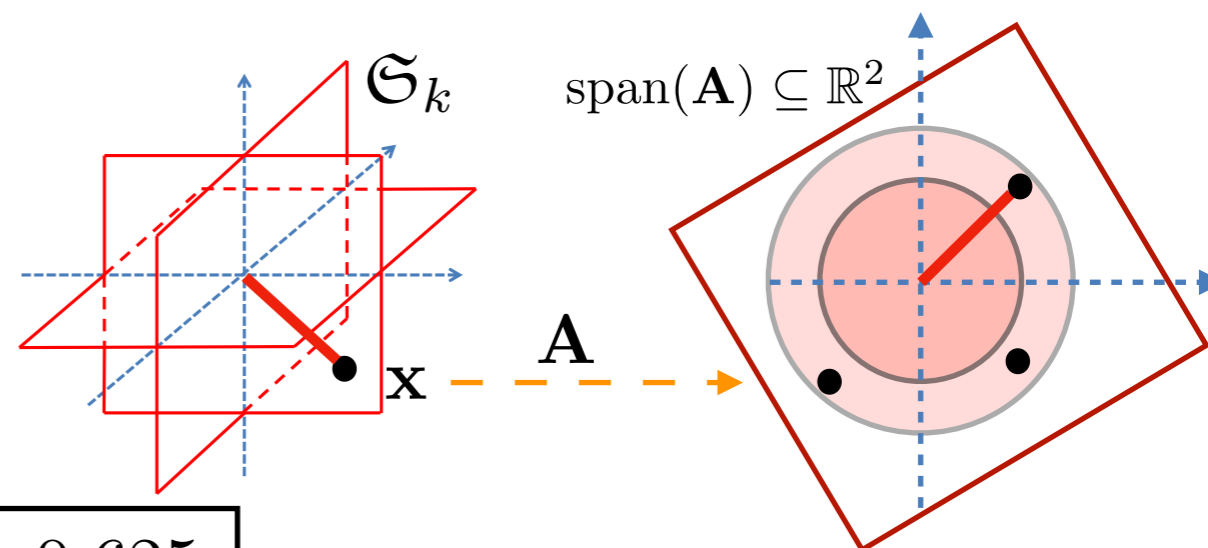
Does the LASSO truly recover $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



The result: Suppose $\delta_{2k} < 4/\sqrt{41} \approx 0.625$

Take $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} = \mathbf{Ax} + \mathbf{e}$
with $\|\mathbf{e}\|_2 \leq \eta$

$\hat{\mathbf{x}}$ sol. of $\begin{cases} \min \|\mathbf{z}\|_1 \\ \text{s.t. } \|\mathbf{y} - \mathbf{Az}\|_2 \leq \eta \end{cases}$

$$\implies \|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \frac{C_1}{\sqrt{k}} d^\circ(\mathbf{x}, \mathfrak{S}_k) + C_2 \eta$$

Compressed sensing theory: an invitation

Guarantees for the LASSO

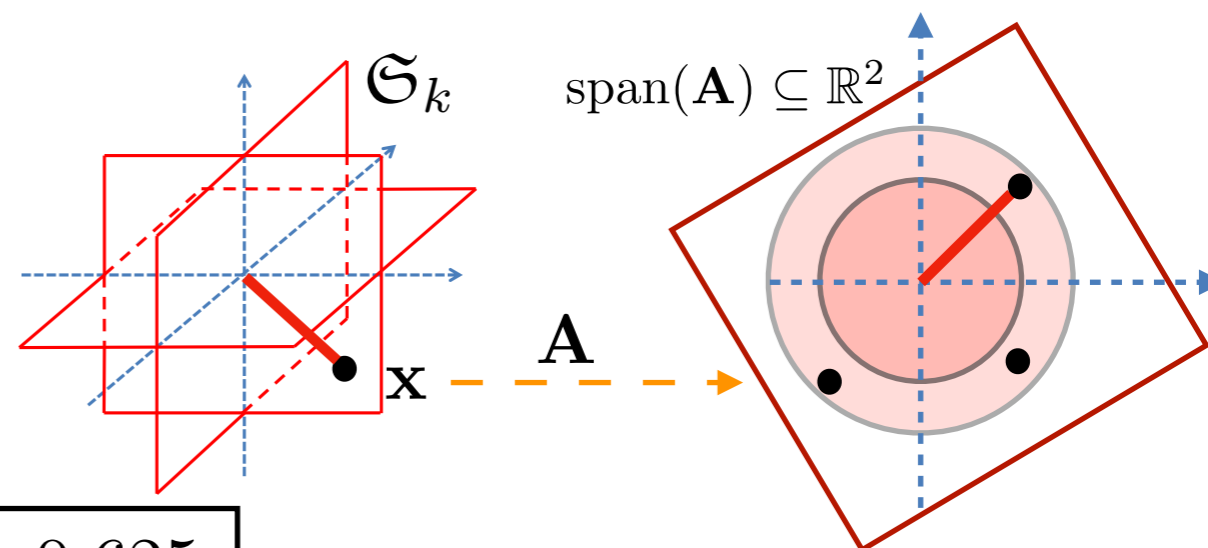
Does the LASSO truly recover $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



The result: Suppose $\delta_{2k} < 4/\sqrt{41} \approx 0.625$

Take $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} = \mathbf{Ax} + \mathbf{e}$
with $\|\mathbf{e}\|_2 \leq \eta$

$\hat{\mathbf{x}}$ sol. of $\begin{cases} \min \|\mathbf{z}\|_1 \\ \text{s.t. } \|\mathbf{y} - \mathbf{Az}\|_2 \leq \eta \end{cases}$

$$\Rightarrow \|\mathbf{x} - \hat{\mathbf{x}}\|_2$$

Error between the estimation and the ground truth

Compressed sensing theory: an invitation

Guarantees for the LASSO

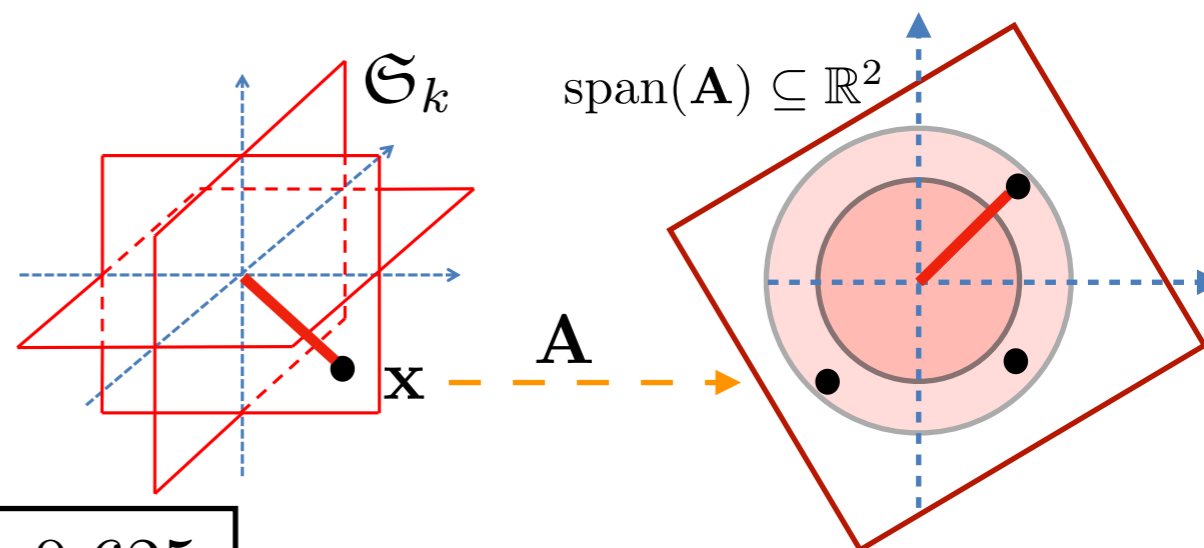
Does the LASSO truly recover $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



The result: Suppose $\delta_{2k} < 4/\sqrt{41} \approx 0.625$

Take $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} = \mathbf{Ax} + \mathbf{e}$
with $\|\mathbf{e}\|_2 \leq \eta$

$\hat{\mathbf{x}}$ sol. of $\begin{cases} \min \|\mathbf{z}\|_1 \\ \text{s.t. } \|\mathbf{y} - \mathbf{Az}\|_2 \leq \eta \end{cases}$

$$\implies \|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \frac{C_1}{\sqrt{k}} d^\circ(\mathbf{x}, \mathcal{S}_k) \dashrightarrow \text{Zero if the vector is exactly } k\text{-sparse}$$

Compressed sensing theory: an invitation

Guarantees for the LASSO

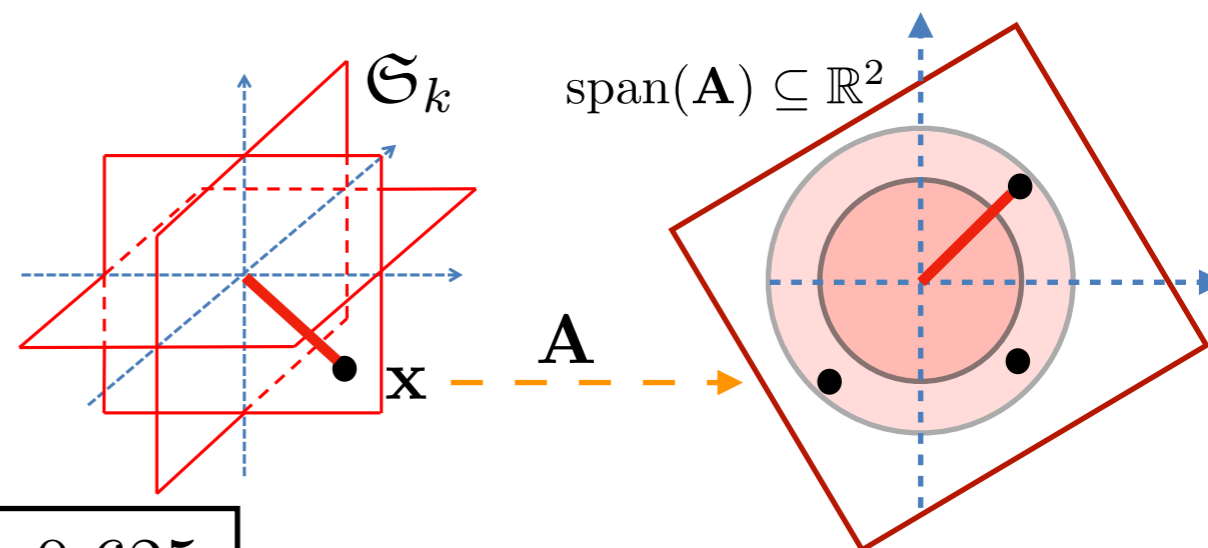
Does the LASSO truly recovers $\mathbf{x} \in \mathbb{R}^d$?



The restricted isometric property (RIP) [Candes & Tao, 2005]

$$\exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse}$$

$$(1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{Ax}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2$$



The result: Suppose $\delta_{2k} < 4/\sqrt{41} \approx 0.625$

Take $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y} = \mathbf{Ax} + \mathbf{e}$
with $\|\mathbf{e}\|_2 \leq \eta$

$\hat{\mathbf{x}}$ sol. of $\begin{cases} \min \|\mathbf{z}\|_1 \\ \text{s.t. } \|\mathbf{y} - \mathbf{Az}\|_2 \leq \eta \end{cases}$

$$\implies \|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \frac{C_1}{\sqrt{k}} d^\circ(\mathbf{x}, \mathfrak{S}_k) + C_2 \eta$$

Zero if there is no noise

Compressed sensing theory: an invitation

How can we recover \mathbf{x} ?

Ill-posed **inverse problem**

Infinity of solutions

Underdetermined $m < d$

$$\mathbf{y} = \mathbf{A} \mathbf{x}$$

The sparse assumption

The true vector lives in a **low-dim** space

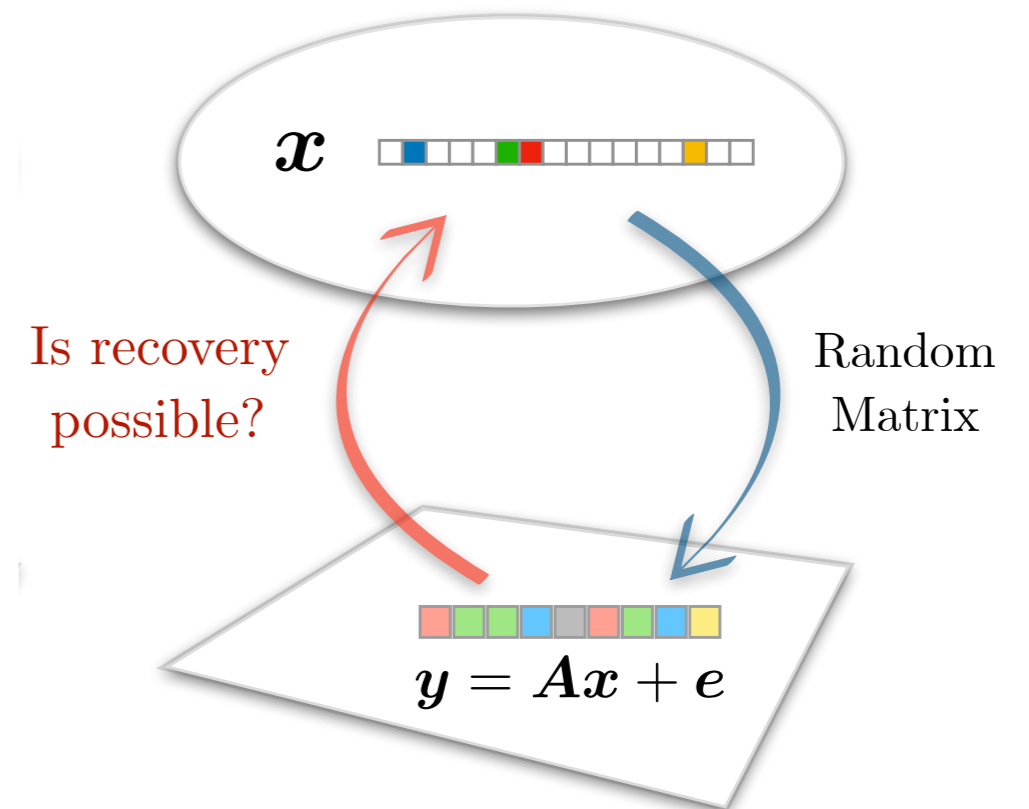
Algorithmic solutions

$$\min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

Theoretical guarantees

E.g. RIP with Gaussian matrices

$$m \gtrsim \delta^{-2} k \ln\left(e \frac{d}{k}\right)$$



| Overview of the talk

- Part I: **A journey in the compressed sensing theory**

- Part II: **A bit of machine learning theory**

- Part III: **The sketching approach**

 - Applied sketching

 - Theoretical guarantees

A bit of machine learning

Fashion Trend Forecasting with AI

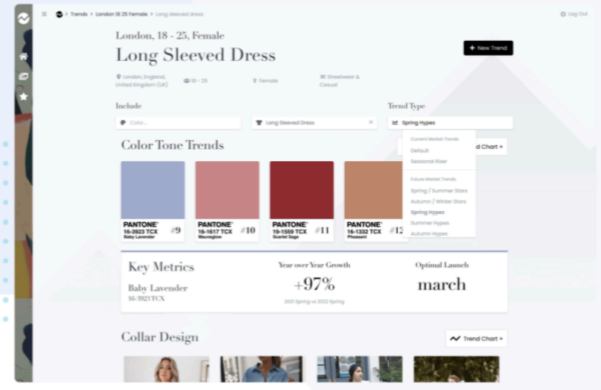

JESSICA MAGALIT - NOVEMBER 9, 2021

🔗 0 🗨️ 0

See What's Next

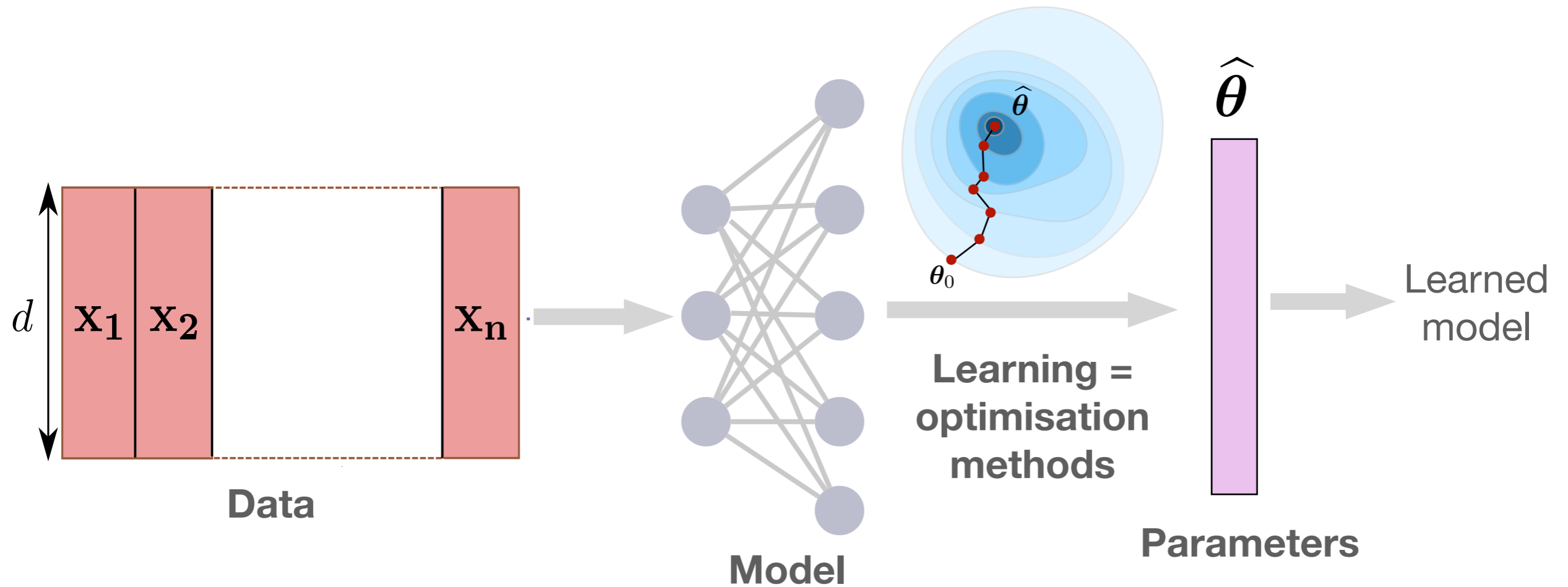
With T-Fashion's AI-powered trend forecasting platform, grasp trend dynamics through billions of interactions taking place online.

Receive customized fashion analytics and data-driven trend insights to produce/buy the right product at the right time.



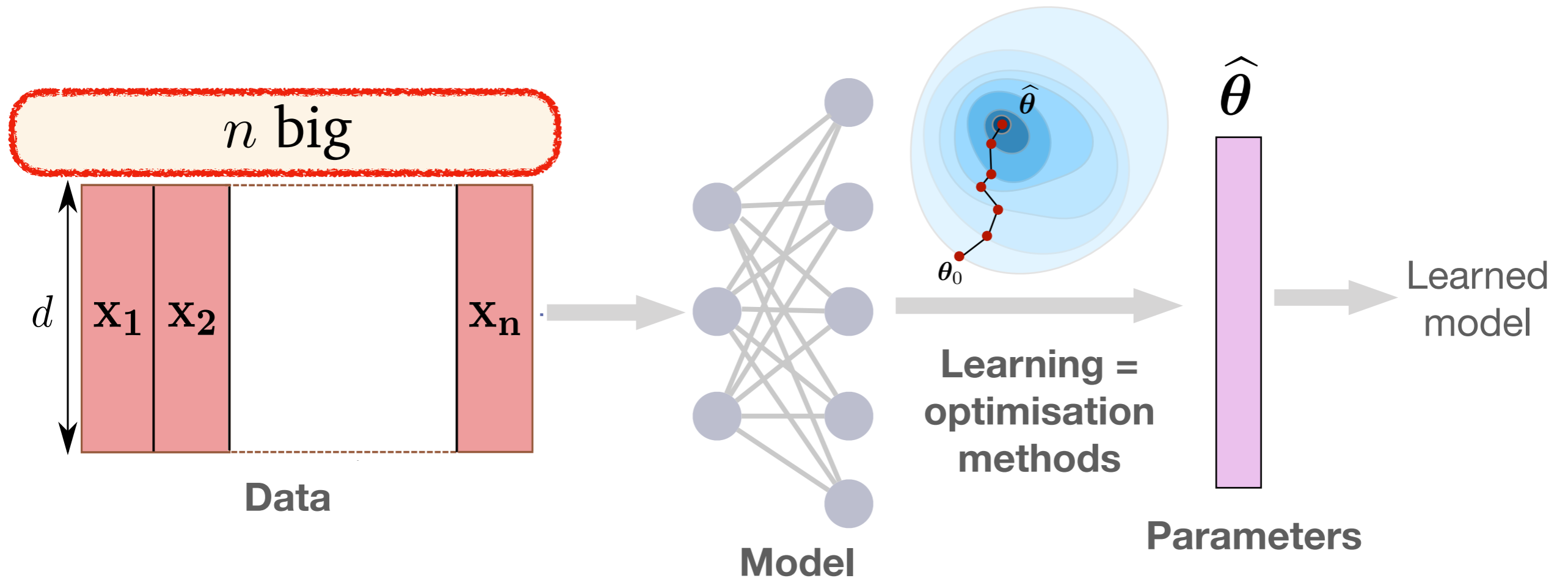
Machine learning theory

■ The big picture:



Machine learning theory

The big picture:

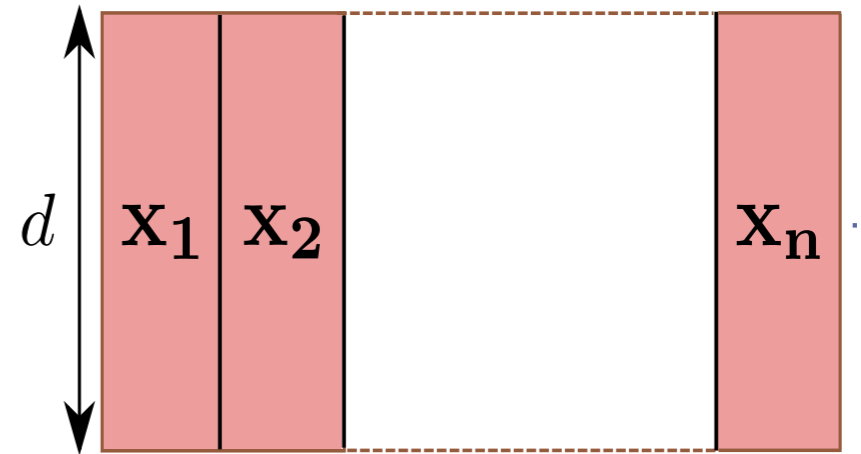


Large scale machine learning

| Machine learning theory

■ The ML setting:

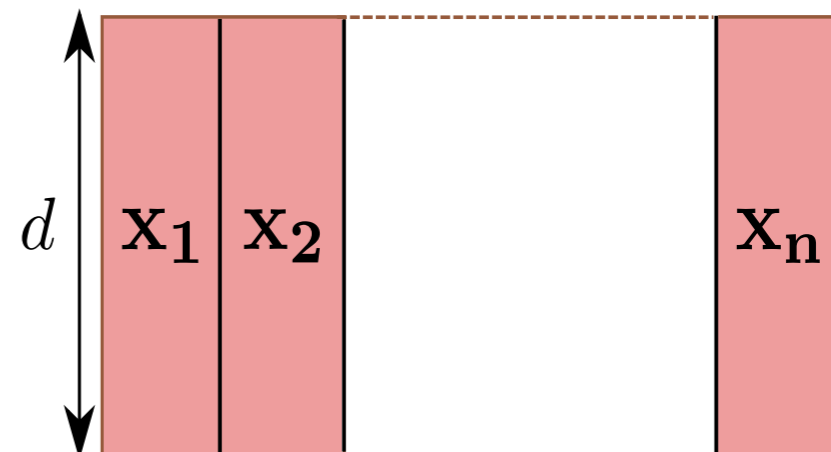
- Data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- Each point $\mathbf{x}_i \sim \pi$
- π is **unknown** and generates the data



Machine learning theory

The ML setting:

- Data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- Each point $\mathbf{x}_i \sim \pi$
- π is **unknown** and generates the data



Empirical risk minimization:

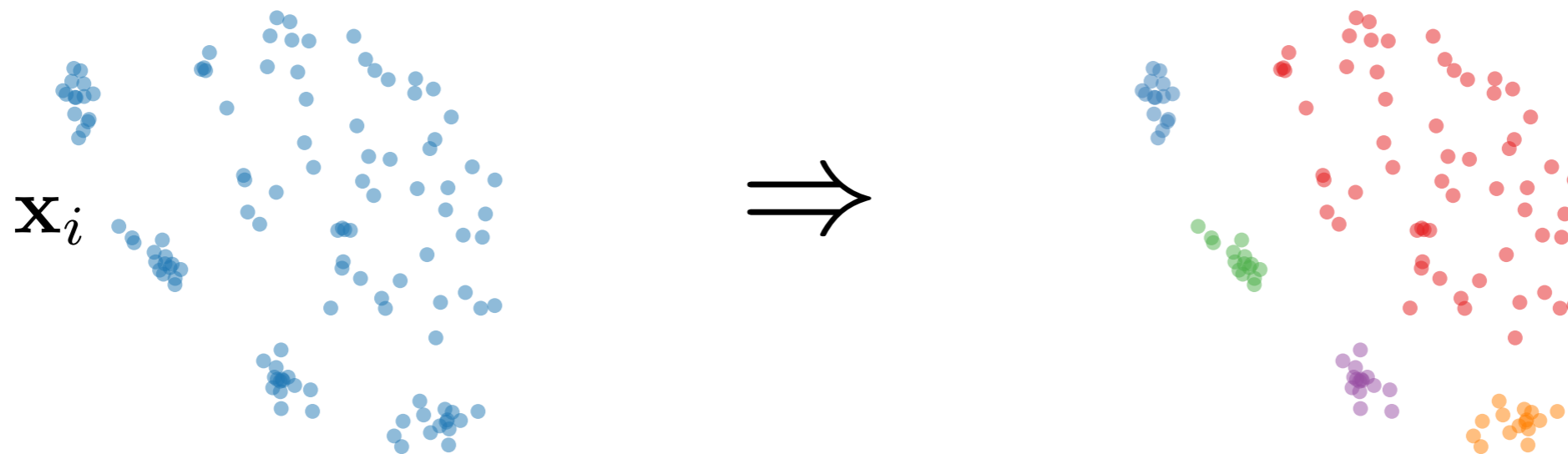
- Find parameters: $\hat{\boldsymbol{\theta}} \in \Theta$
- That minimizes:

$$\min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \boldsymbol{\theta}) + \lambda \text{Reg}(\boldsymbol{\theta})$$

- Where $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}$ is a **loss** and Reg a **regularization** term

Machine learning theory

Unsupervised learning: K-means clustering



Organize training samples **in groups** (say K)

Find K **clusters** that **best** represent our data

We look for $\theta = (\mathbf{c}_1, \dots, \mathbf{c}_K)$, $\mathbf{c}_k \in \mathbb{R}^d$

The loss is $\ell(\mathbf{x}_i, \theta) = \min_{k \in \llbracket K \rrbracket} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$

↓
Squared distance between the point and its closest cluster

Machine learning theory

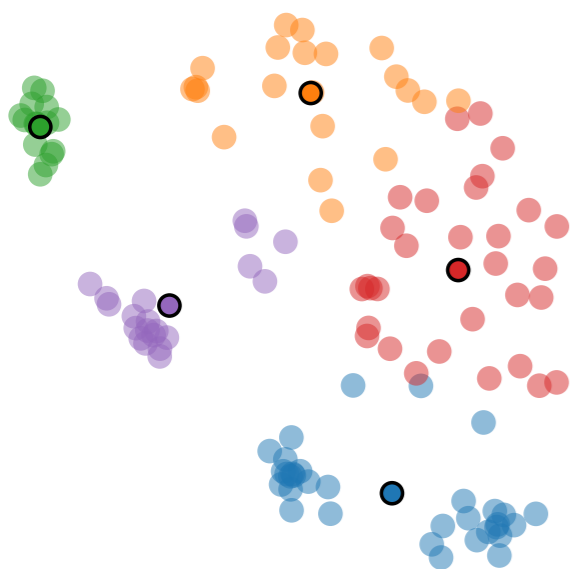


Unsupervised learning: K-means clustering

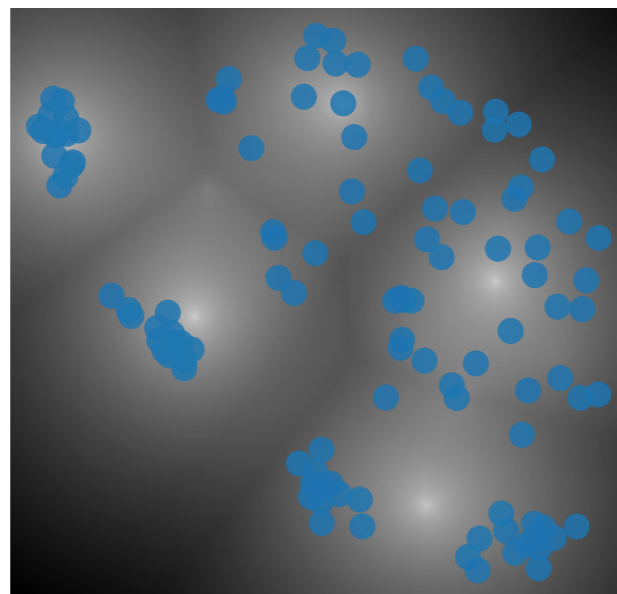
We aim at solving: [Steinhaus & al, 1956, McQuenn & al, 1967]

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \frac{1}{n} \sum_{i=1}^n \min_{k \in \llbracket K \rrbracket} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

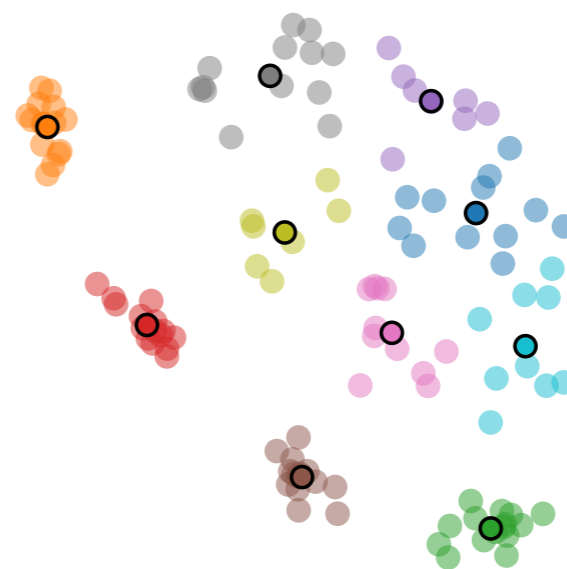
K-means for K=5



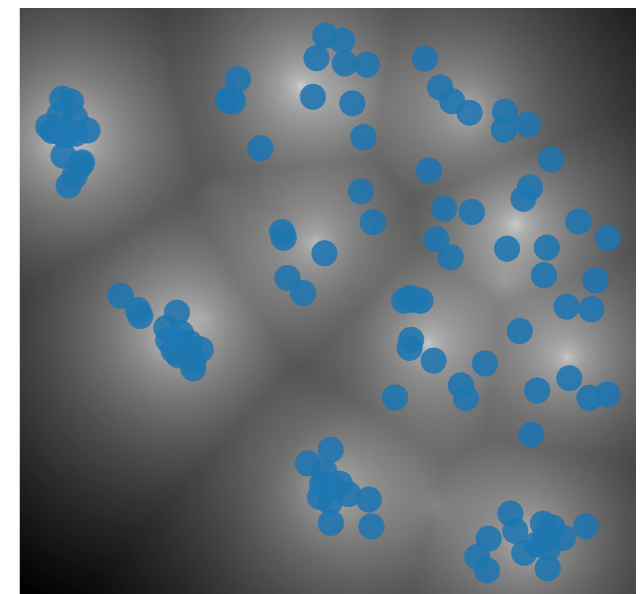
$f(x) = \min_k |x - c_k|^2$ for K=5



K-means for K=10



$f(x) = \min_k |x - c_k|^2$ for K=10

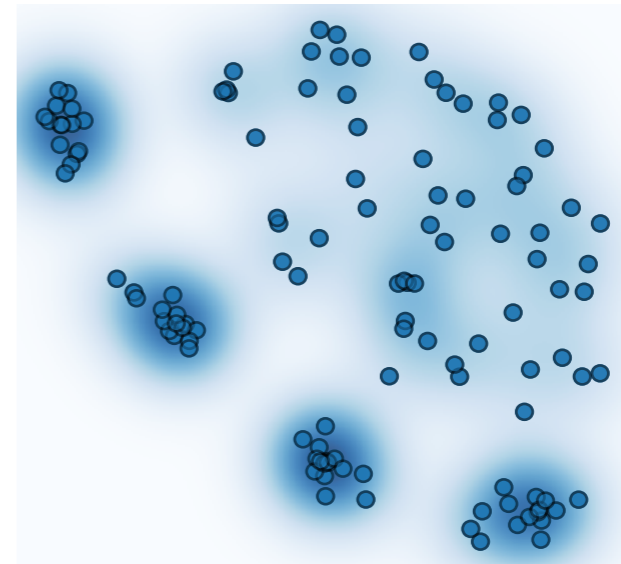
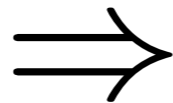
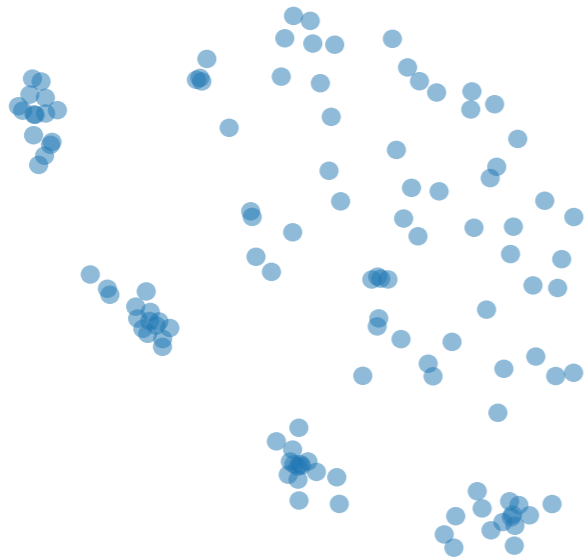


It is a NP-Hard problem: can be tackled by **Lloyd's algorithm**

Complexity (in time): $\mathcal{O}(nKd)$

| Machine learning theory

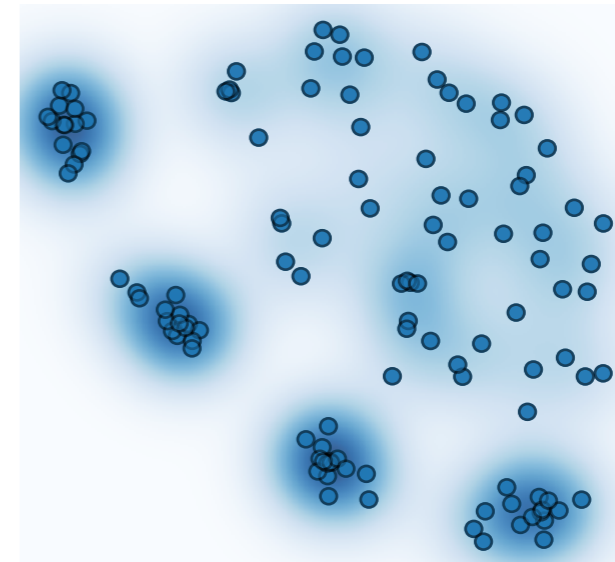
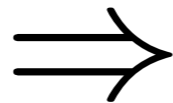
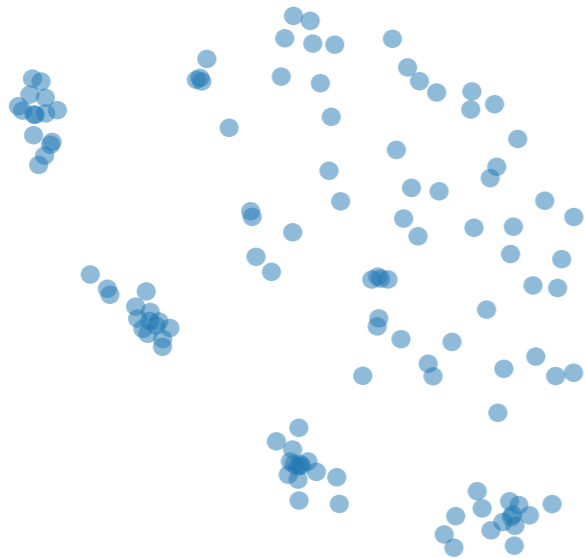
■ Unsupervised learning: GMM



- Estimate a **probability density** $\hat{\pi}$ from the **samples**
 - $\forall \mathbf{x}, \hat{\pi}(\mathbf{x}) \geq 0, \int \hat{\pi}(\mathbf{x}) d\mathbf{x} = 1$
- } Density estimation

| Machine learning theory

■ Unsupervised learning: **GMM**



- Estimate a **probability density** $\hat{\pi}$ from the **samples**
 - $\forall \mathbf{x}, \hat{\pi}(\mathbf{x}) \geq 0, \int \hat{\pi}(\mathbf{x}) d\mathbf{x} = 1$
- } Density estimation

■ Parametrized probability distributions

- Look only for a family of distributions given by param. θ
- E.g. **Gaussian** $\theta = \{\mu, \Sigma\}$

$$\pi_{\mu, \Sigma}(\mathbf{x}) := (2\pi)^{-\frac{d}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)\right)$$

| Machine learning theory

■ Unsupervised learning: GMM

- The model is a mixture of Gaussian

$$\theta = \{\alpha_k, \mu_k, \Sigma_k\}_{k=1}^K$$

$$\pi_{\theta}(\mathbf{x}) = \sum_{k=1}^K \alpha_k \pi_{\mu_k, \Sigma_k}(\mathbf{x})$$

Machine learning theory

Unsupervised learning: GMM

- The model is a **mixture of Gaussian**

$$\theta = \{\alpha_k, \mu_k, \Sigma_k\}_{k=1}^K$$

$$\pi_{\theta}(\mathbf{x}) = \sum_{k=1}^K \alpha_k \pi_{\mu_k, \Sigma_k}(\mathbf{x})$$

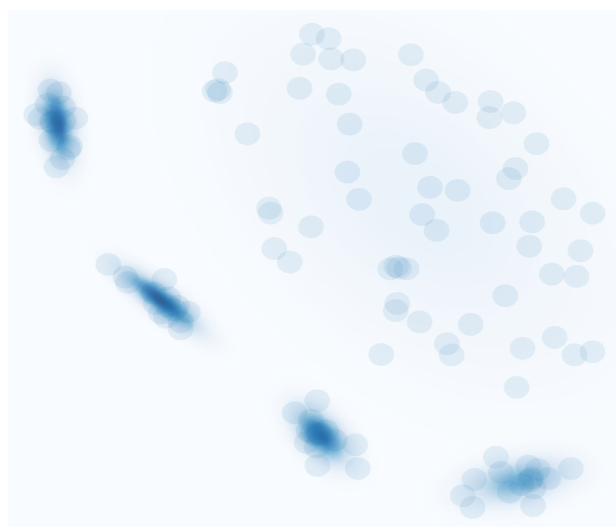
- We aim at solving **the MLE problem**: [Dempster & al, 1977]

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n -\log(\pi_{\theta}(\mathbf{x}_i))$$

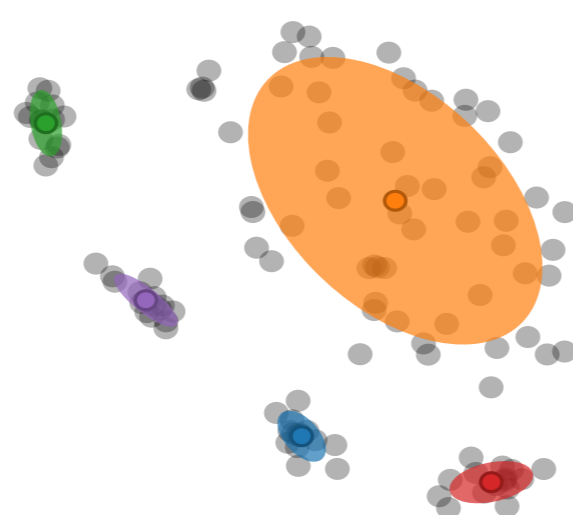


- Here $\ell(\mathbf{x}_i, \theta) = -\log(\pi_{\theta}(\mathbf{x}_i))$ (neg log likelihood)

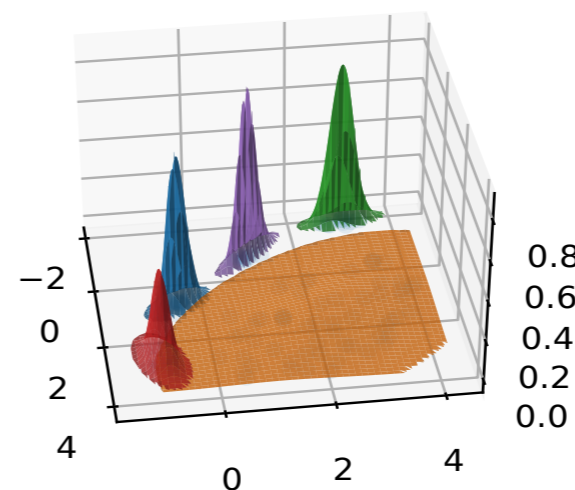
GMM density



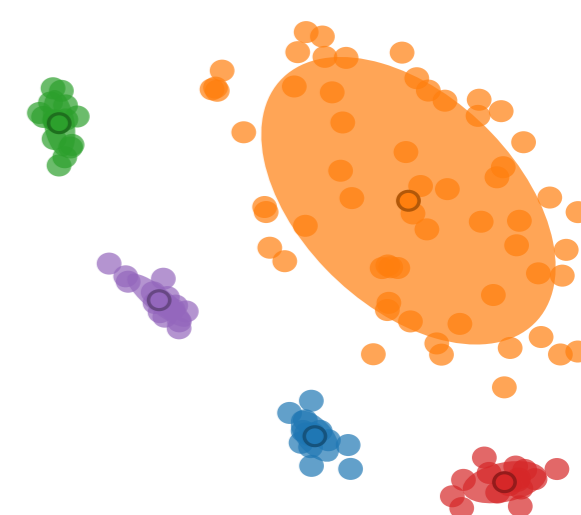
Estimated GMM



GMM mixture densities



GMM clustering



| Machine learning theory

■ Supervised learning:

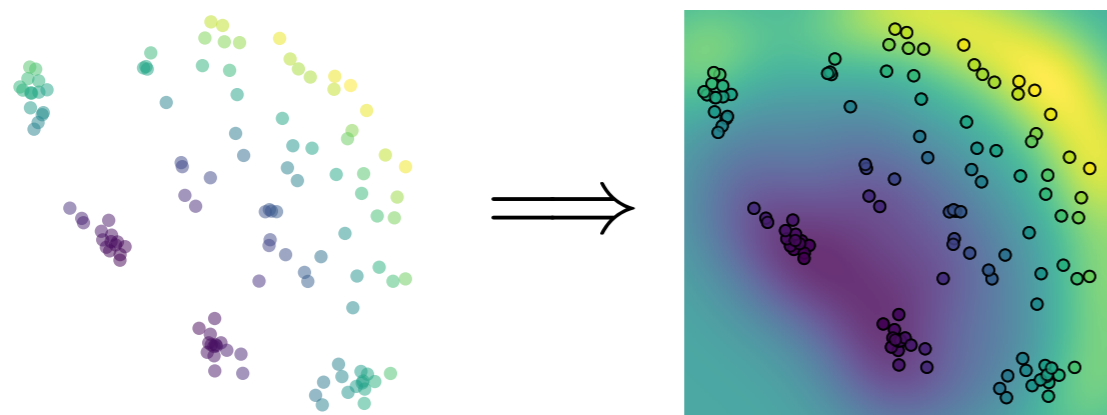
■ In the supervised setting: $\mathbf{x}_i = (\mathbf{z}_i, y_i)$

Machine learning theory

Supervised learning:

In the supervised setting: $\mathbf{x}_i = (\mathbf{z}_i, y_i)$

Regression: $y_i \in \mathbb{R}$ $h(\mathbf{z}_i) \in \mathbb{R}$



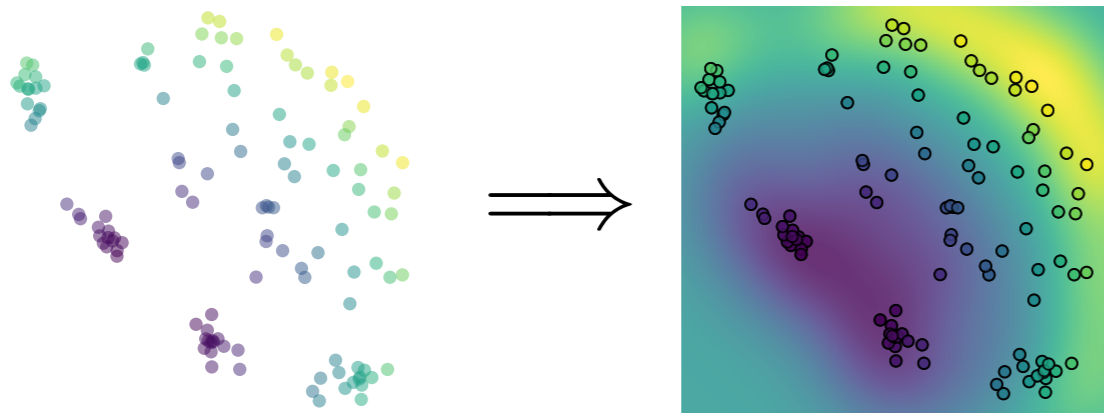
$$\ell(\mathbf{x}_i, h) = \|y_i - h(\mathbf{z}_i)\|_2^2$$

Machine learning theory

Supervised learning:

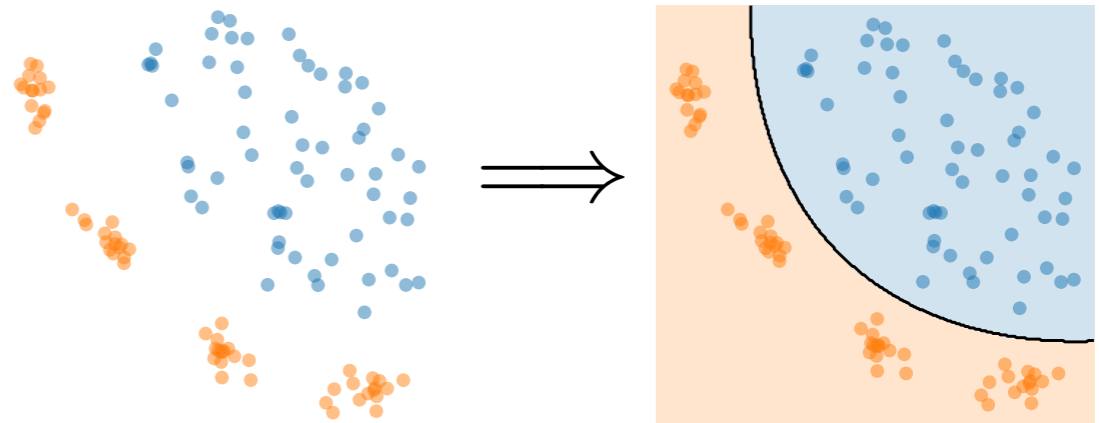
In the supervised setting: $\mathbf{x}_i = (\mathbf{z}_i, y_i)$

Regression: $y_i \in \mathbb{R}$ $h(\mathbf{z}_i) \in \mathbb{R}$



$$\ell(\mathbf{x}_i, h) = \|y_i - h(\mathbf{z}_i)\|_2^2$$

Classification: $y_i, h(\mathbf{z}_i) \in \{-1, 1\}$

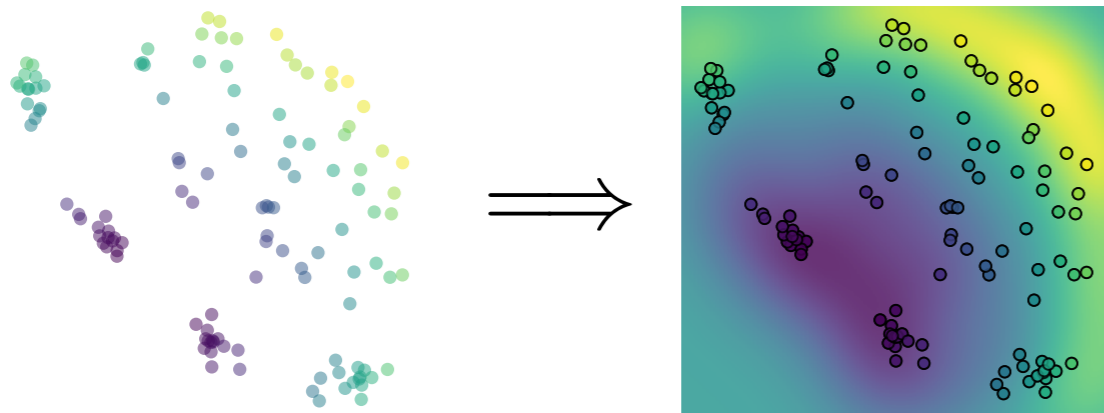


Machine learning theory

Supervised learning:

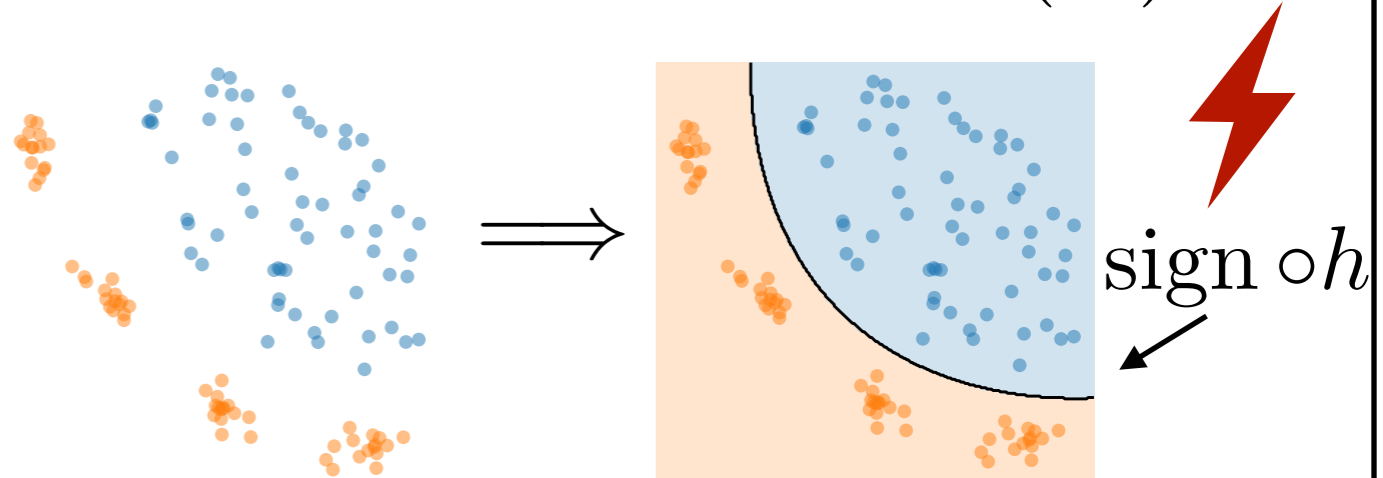
In the supervised setting: $\mathbf{x}_i = (\mathbf{z}_i, y_i)$

Regression: $y_i \in \mathbb{R}$ $h(\mathbf{z}_i) \in \mathbb{R}$



$$\ell(\mathbf{x}_i, h) = \|y_i - h(\mathbf{z}_i)\|_2^2$$

Classification: $y_i \in \{-1, 1\}$ $h(\mathbf{z}_i) \in \mathbb{R}$

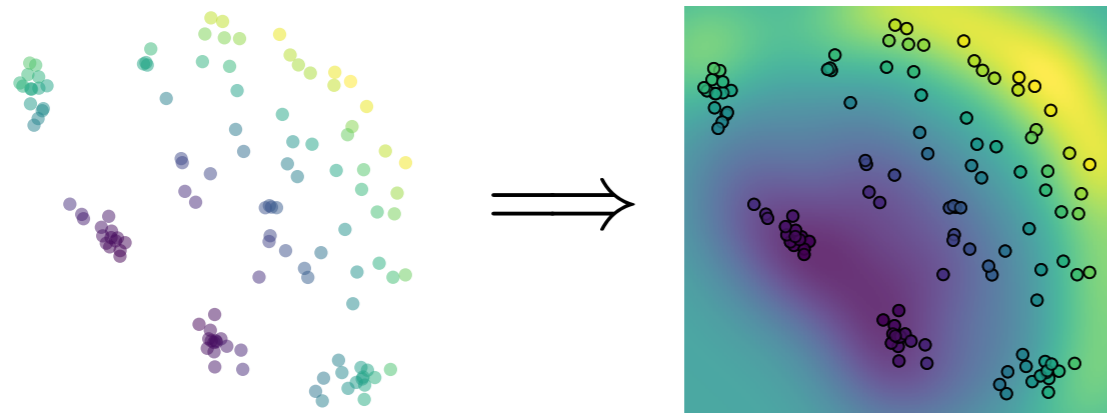


Machine learning theory

Supervised learning:

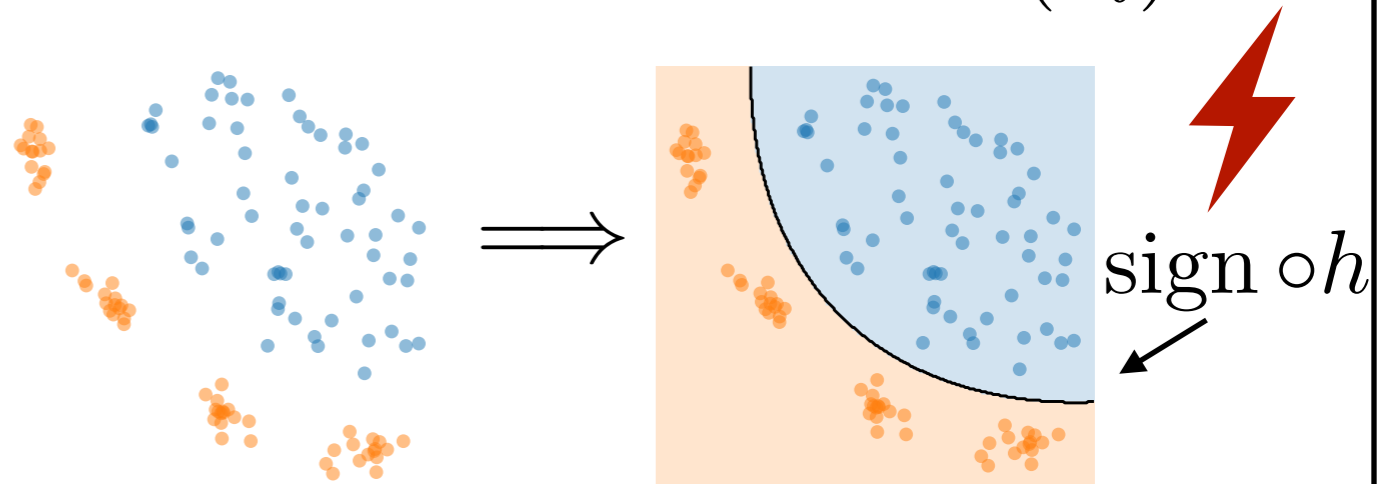
In the supervised setting: $\mathbf{x}_i = (\mathbf{z}_i, y_i)$

Regression: $y_i \in \mathbb{R}$ $h(\mathbf{z}_i) \in \mathbb{R}$



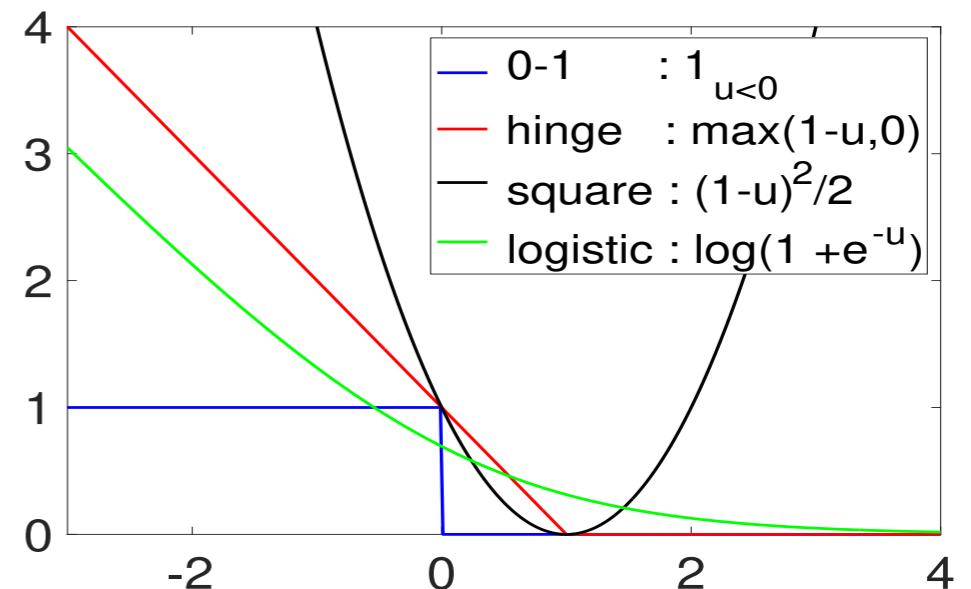
$$\ell(\mathbf{x}_i, h) = \|y_i - h(\mathbf{z}_i)\|_2^2$$

Classification: $y_i \in \{-1, 1\}$ $h(\mathbf{z}_i) \in \mathbb{R}$



$$\ell(\mathbf{x}_i, h) = \phi(y_i h(\mathbf{x}_i))$$

$\phi(u)$

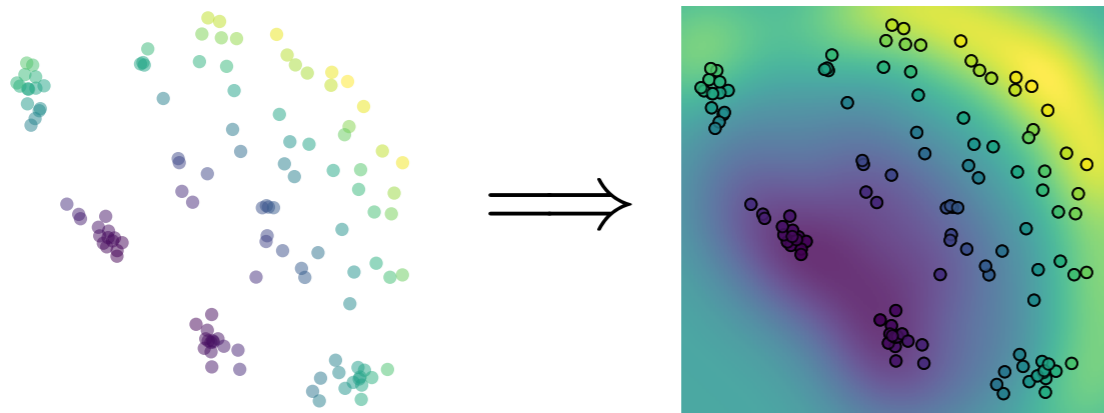


Machine learning theory

Supervised learning:

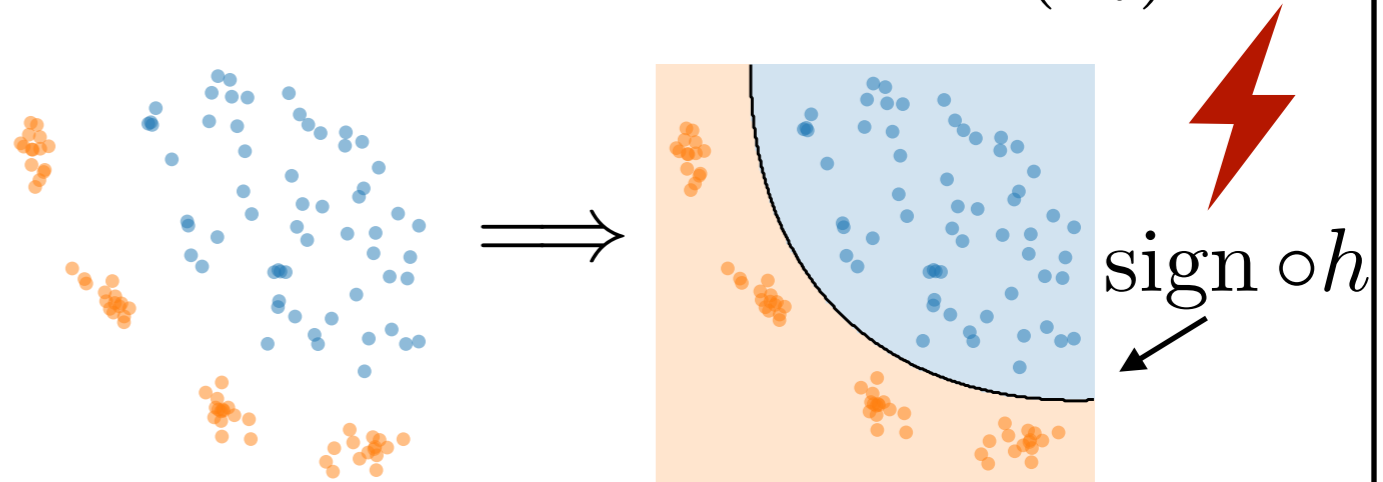
In the supervised setting: $\mathbf{x}_i = (\mathbf{z}_i, y_i)$

Regression: $y_i \in \mathbb{R}$ $h(\mathbf{z}_i) \in \mathbb{R}$

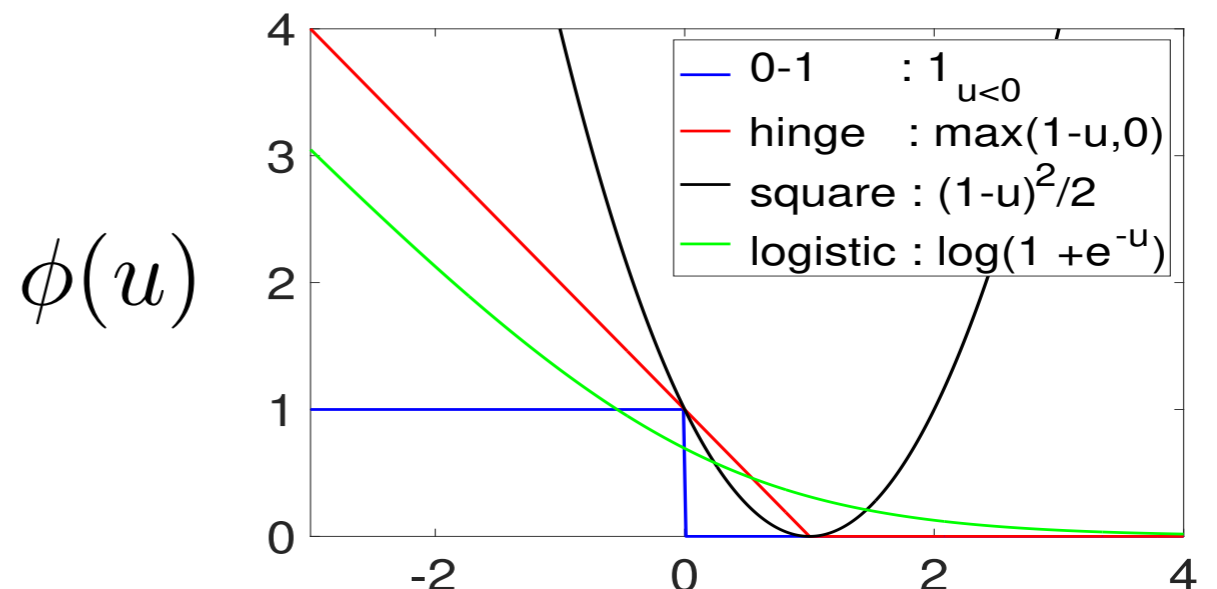


$$\ell(\mathbf{x}_i, h) = \|y_i - h(\mathbf{z}_i)\|_2^2$$

Classification: $y_i \in \{-1, 1\}$ $h(\mathbf{z}_i) \in \mathbb{R}$



$$\ell(\mathbf{x}_i, h) = \phi(y_i h(\mathbf{x}_i))$$

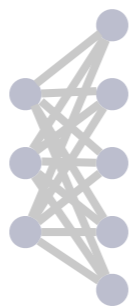


Parametrized model

$$h = f_{\theta}$$

Linear reg: $f_{\theta}(\mathbf{z}) = \theta^{\top} \mathbf{z}$

Neural networks: $f_{\theta}(\mathbf{z}) = \text{NN}_{\theta}(\mathbf{z})$



| Machine learning theory

■ ML in practice

■ ERM

$$\min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \boldsymbol{\theta}) + \lambda \text{Reg}(\boldsymbol{\theta})$$

Machine learning theory

ML in practice

ERM

$$\min_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \boldsymbol{\theta}) + \lambda \text{Reg}(\boldsymbol{\theta})$$

everything differentiable

Machine learning theory

ML in practice:

ERM

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ - - - - \blacktriangleright you **want** to minimize it



Machine learning theory

ML in practice:

ERM

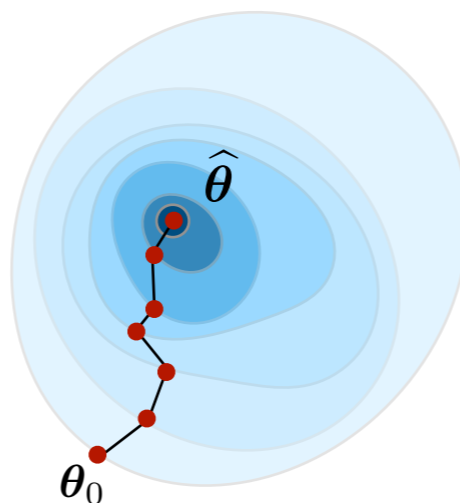
$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ - - - - \blacktriangleright you want to minimize it



Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{R}_n(\theta)$$



Machine learning theory

ML in practice:

ERM

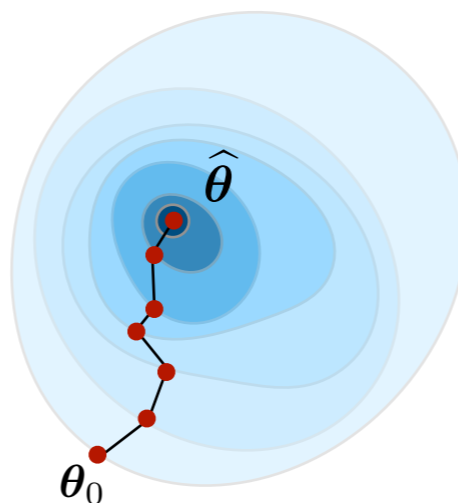
$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ - - - - \blacktriangleright you want to minimize it



Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{R}_n(\theta)$$



We are lazy:  PyTorch

$$\theta_{k+1} = \theta_k - \eta \text{autodiff}[\mathcal{R}_n(\theta)]$$

Machine learning theory

ML in practice:

ERM

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ - - - - \blacktriangleright you want to minimize it



Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{R}_n(\theta)$$

Many many many variants

Fix step-size, or change it $(\eta_k)_k$

Momentum, averaging, adaptative step-size strategies:

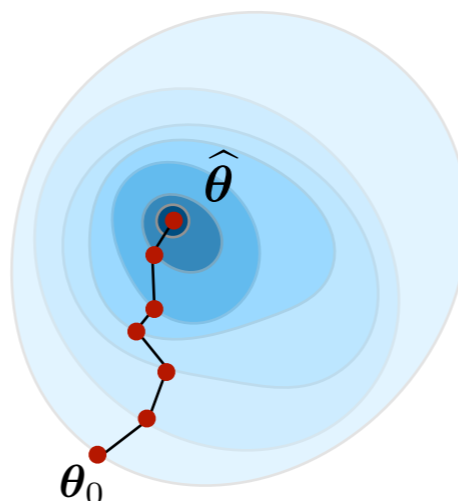
Momentum and Accelerated gradients [Nesterov, 1983]

RMSPROP [Tieleman & Hinton, 2012]

Adam [Kingma & Ba, 2014]

We are lazy:  PyTorch

$$\theta_{k+1} = \theta_k - \eta \text{autodiff}[\mathcal{R}_n(\theta)]$$



Machine learning theory

ML in practice:

ERM

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ - - - - \blacktriangleright you want to minimize it



Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{R}_n(\theta)$$

Many many many variants

Fix step-size, or change it $(\eta_k)_k$

Momentum, averaging, adaptative step-size strategies:

Momentum and Accelerated gradients [Nesterov, 1983]

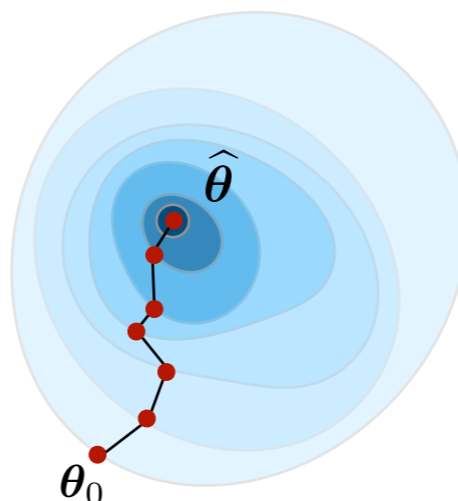
RMSPROP [Tieleman & Hinton, 2012]

Adam [Kingma & Ba, 2014]

We are lazy:  PyTorch

$$\theta_{k+1} = \theta_k - \eta \text{autodiff}[\mathcal{R}_n(\theta)]$$

Hope that $\hat{\theta} = \theta_{\infty}$ is « good »



| Machine learning theory

Hope that $\hat{\theta} = \theta_{\infty}$ is « good » ??

Machine learning theory

Hope that $\hat{\theta} = \theta_{\infty}$ is « good » ??

Approximate the true risk

ERM: $\hat{\theta} \in \arg \min_{\theta \in \Theta} \mathcal{R}_n(\theta)$

$$= \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, f_{\theta})$$

Machine learning theory

Hope that $\hat{\theta} = \theta_{\infty}$ is « good » ??

Approximate the true risk

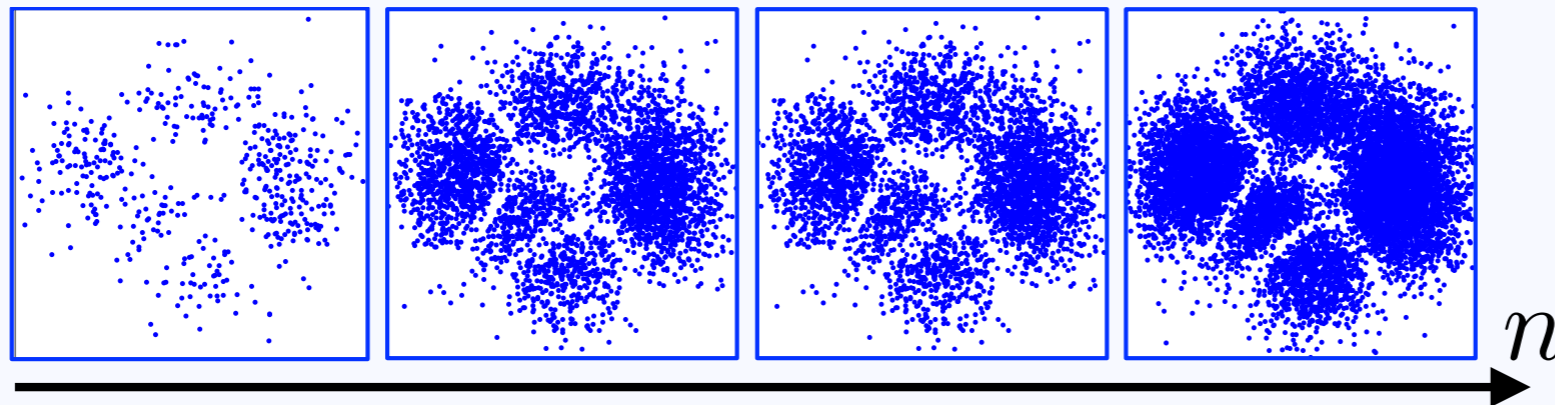
ERM: $\hat{\theta} \in \arg \min_{\theta \in \Theta} \mathcal{R}_n(\theta)$

$\mathcal{R}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta)$

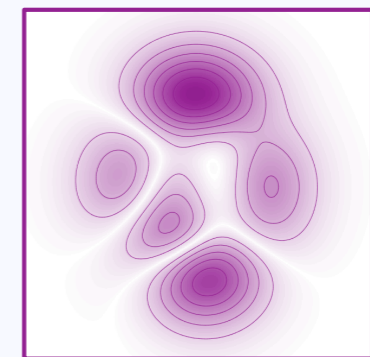
$\mathcal{R}(\theta) = \mathbb{E}_{\mathbf{x} \sim \pi} [\ell(\mathbf{x}, \theta)]$

Side note

Empirical distribution $\pi_n = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$



True distrib. π



$n = \infty$

Machine learning theory

Hope that $\hat{\theta} = \theta_{\infty}$ is « good » ??

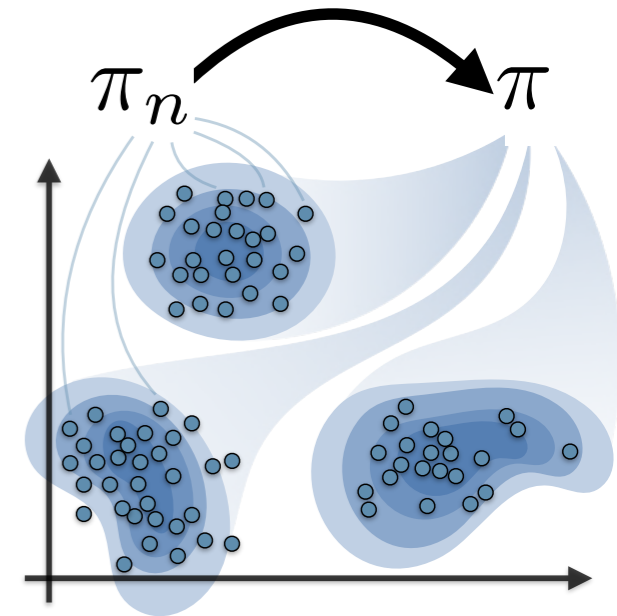
Approximate the true risk

■ **ERM:** $\hat{\theta} \in \arg \min_{\theta \in \Theta} \mathcal{R}_n(\theta) = \mathbb{E}_{\mathbf{x} \sim \pi_n} [\ell(\mathbf{x}, \theta)]$

■ **True risk:** $\mathcal{R}(\theta) = \mathbb{E}_{\mathbf{x} \sim \pi} [\ell(\mathbf{x}, \theta)]$

■ Best param. (if we have all the data in the world)

$$\theta^* \in \arg \min_{\theta \in \Theta} \mathcal{R}(\theta)$$



Machine learning theory

Hope that $\hat{\theta} = \theta_{\infty}$ is « good » ??

Approximate the true risk

■ **ERM:** $\hat{\theta} \in \arg \min_{\theta \in \Theta} \mathcal{R}_n(\theta) = \mathbb{E}_{\mathbf{x} \sim \pi_n} [\ell(\mathbf{x}, \theta)]$

■ **True risk:** $\mathcal{R}(\theta) = \mathbb{E}_{\mathbf{x} \sim \pi} [\ell(\mathbf{x}, \theta)]$

■ Best param. (if we have all the data in the world)

$$\theta^* \in \arg \min_{\theta \in \Theta} \mathcal{R}(\theta)$$

Wished guarantees:

$$0 \leq \mathcal{R}(\hat{\theta}) - \mathcal{R}(\theta^*) \leq \lambda_n \xrightarrow[n \rightarrow +\infty]{} 0$$

Machine learning theory

ML in practice:

ERM

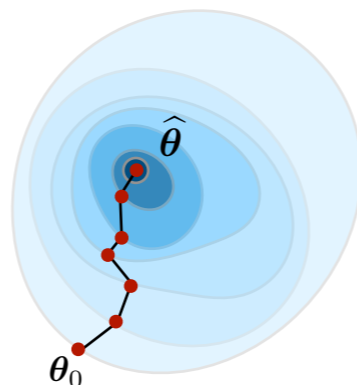
$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ — — — — \blacktriangleright you **want** to minimize it



Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{R}_n(\theta)$$



Guarantees:

Hope that $\hat{\theta} = \theta_{\infty}$ is « **good** »

Machine learning theory

ML in practice:

ERM

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ — — — — \blacktriangleright you **want** to minimize it



Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{R}_n(\theta)$$

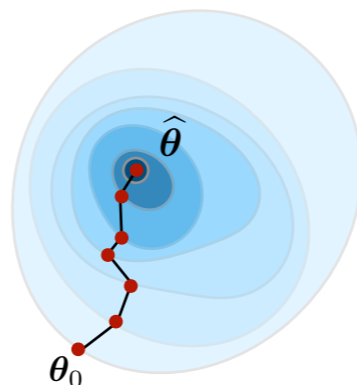
$$\downarrow \sum_{i=1}^n \nabla_{\theta} \ell(\mathbf{x}_i, \theta)$$

Expensive when large scale

Complexity: $\mathcal{O}(n \times C_{\nabla} \times n_{it})$ may not even fit in memory

Alternative: SGD, mini-batches [Bottou, 2010]

BUT requires multiple passes (epochs)



Guarantees:

Hope that $\hat{\theta} = \theta_{\infty}$ is « **good** »

Machine learning theory

ML in practice:

ERM

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$

Empirical risk: $\mathcal{R}_n(\theta)$ — — — — \blacktriangleright you **want** to minimize it



Gradient descent:

$$\theta_{k+1} = \theta_k - \eta \nabla_{\theta} \mathcal{R}_n(\theta)$$

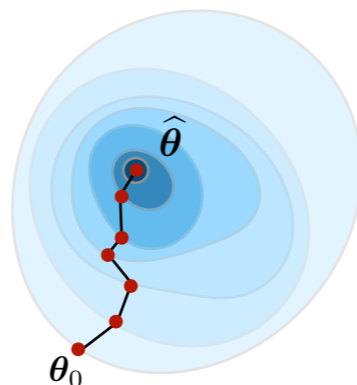
$$\downarrow \sum_{i=1}^n \nabla_{\theta} \ell(\mathbf{x}_i, \theta)$$

Expensive when large scale

Complexity: $\mathcal{O}(n \times C_{\nabla} \times n_{it})$

Alternative: SGD, mini-batches [Bottou, 2010]

BUT requires multiple passes (epochs)



Guarantees:

Hope that $\hat{\theta} = \theta_{\infty}$ is « **good** »

Difficult to prove

Non-convexity/ high-dim

Stats/Optim/Approx theory

Machine learning theory

ML in practice:

ERM

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \theta) + \lambda \text{Reg}(\theta)$$



Empirical risk $\mathcal{P}_n(\theta)$

Gradient

Compressive learning theory

$$\theta_{k+1} = \theta_k - \eta \sum_{i=1}^n \nabla_{\theta} \ell(\mathbf{x}_i, \theta)$$

« good »

$$\sum_{i=1}^n \nabla_{\theta} \ell(\mathbf{x}_i, \theta)$$

Expensive when large scale

Complexity: $\mathcal{O}(n \times C_{\nabla} \times n_{it})$

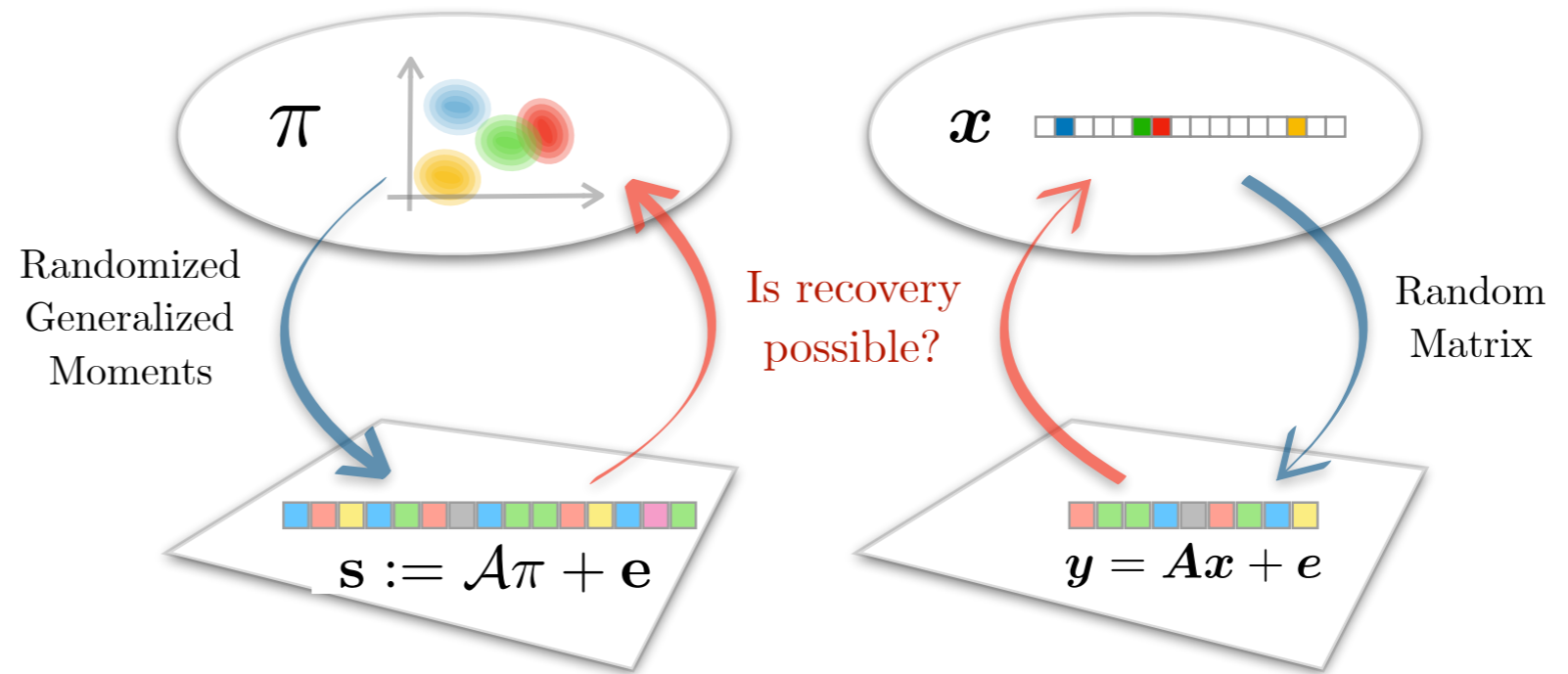
Alternative: SGD, mini-batches [Bottou, 2010]

BUT requires multiple passes (epochs)

Difficult to prove

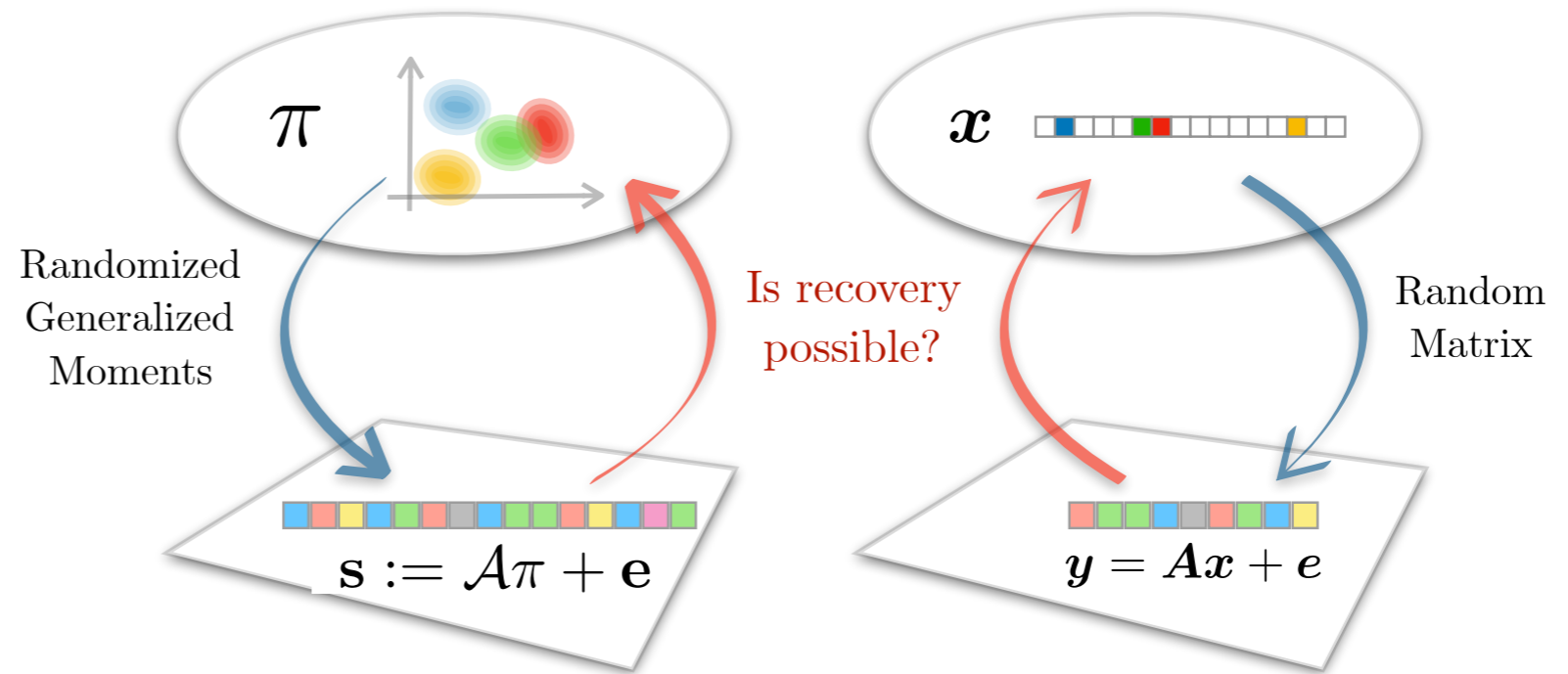
Non-convexity/ high-dim

Stats/Optim/Approx theory



Compressive Learning

- Theory of sketching
- Sketching in practice
- Theoretical guarantees
- Limitations & perspectives

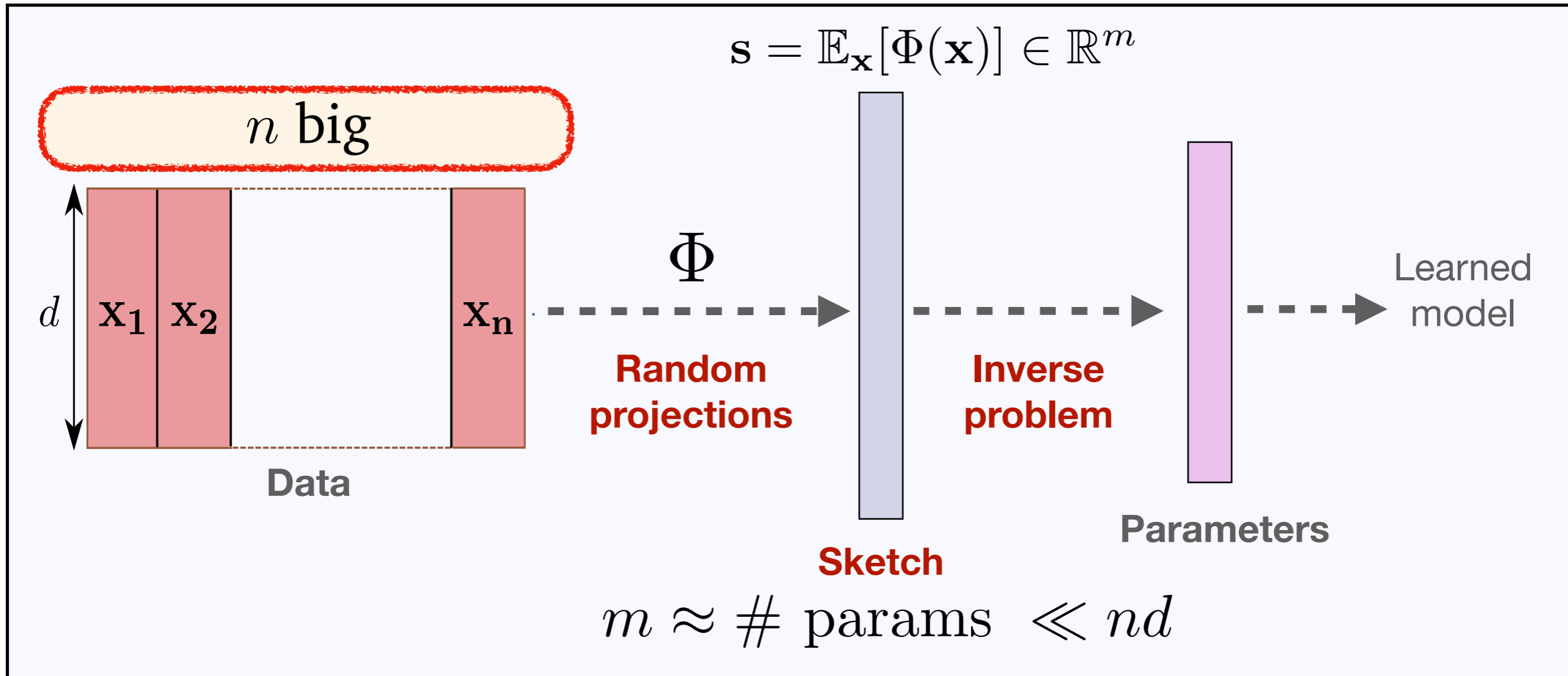


Compressive Learning

- Theory of sketching
- Sketching in practice
- Theoretical guarantees
- Limitations & perspectives

Theory of sketching

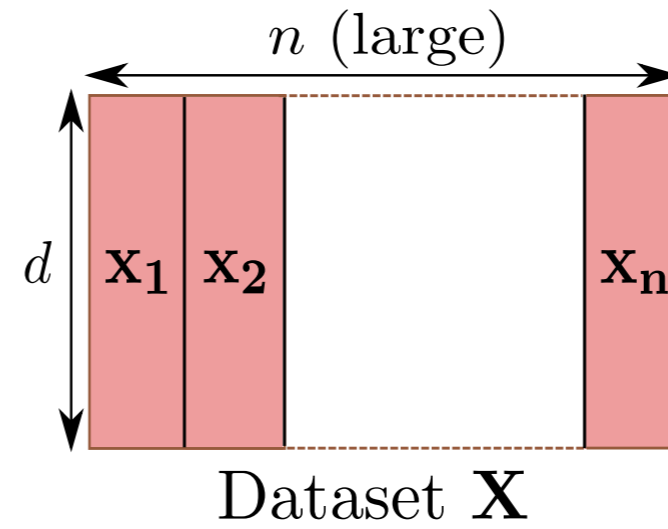
The big picture



| Theory of sketching

■ « Dimension » reduction

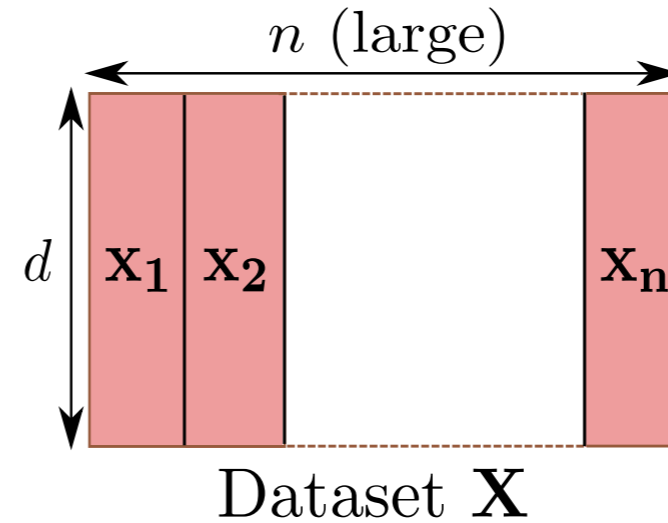
■ « Low-dim » representation
of a dataset



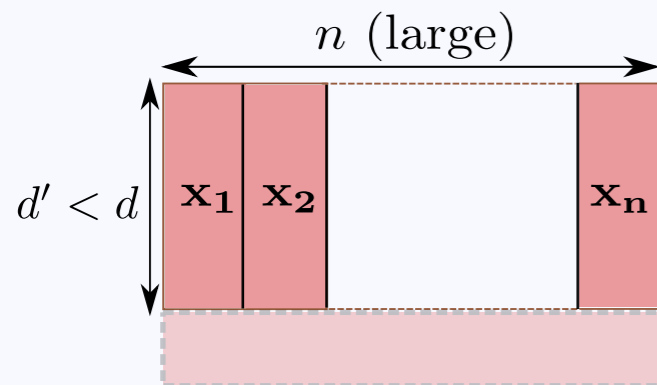
Theory of sketching

« Dimension » reduction

« Low-dim » representation of a dataset



Dimension reduction

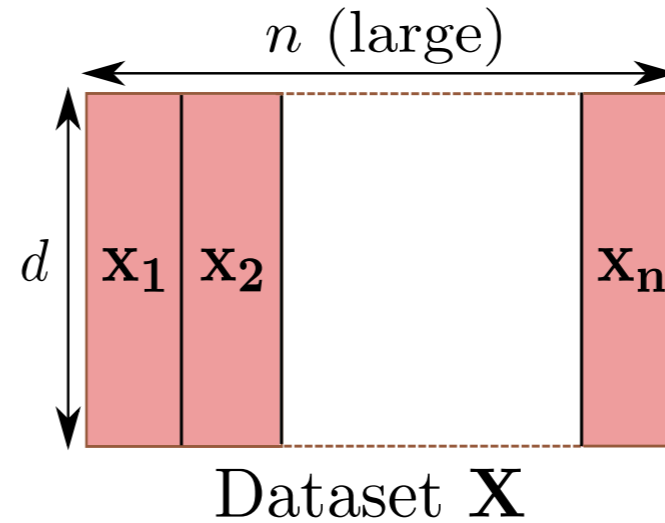


- Random projections (JL lemma)
- Feature selection
- Minimum distortion embedding, PCA

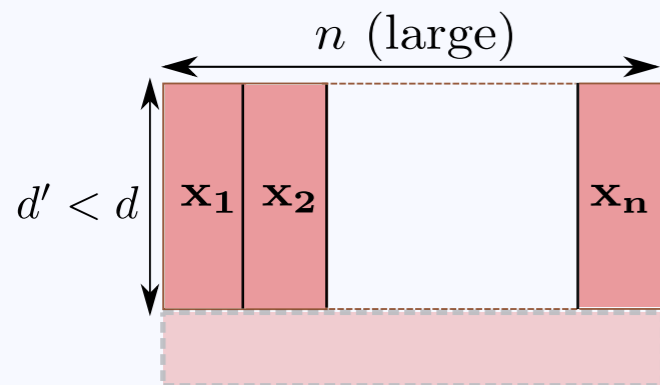
Theory of sketching

« Dimension » reduction

« Low-dim » representation of a dataset

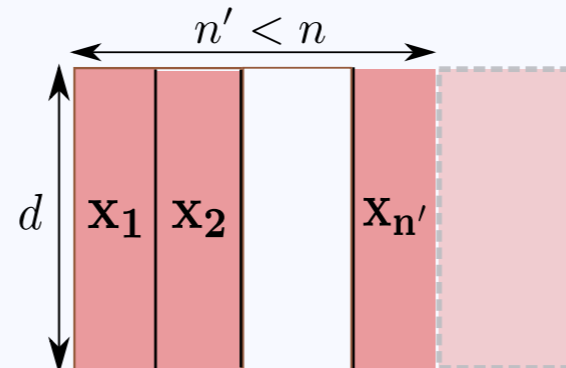


Dimension reduction



- | Random projections (JL lemma)
- | Feature selection
- | Minimum distortion embedding, PCA

Subsampling

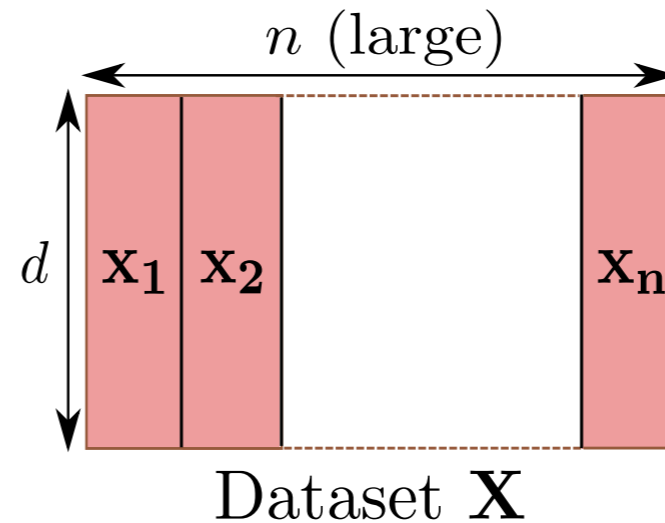


- | Coresets
- | Importance sampling

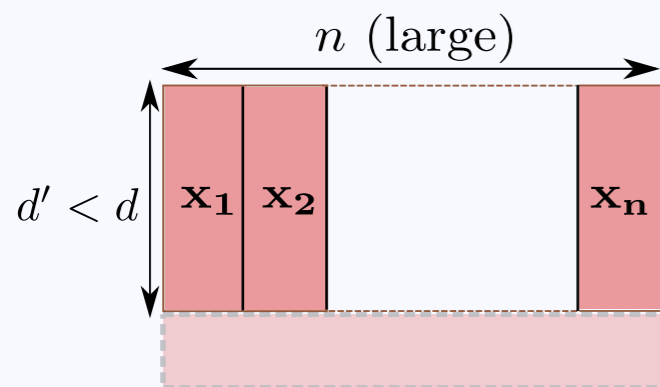
Theory of sketching

« Dimension » reduction

« Low-dim » representation of a dataset



Dimension reduction

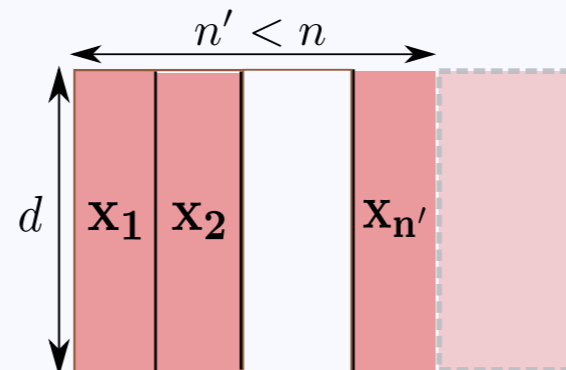


Random projections (JL lemma)

Feature selection

Minimum distortion embedding, PCA

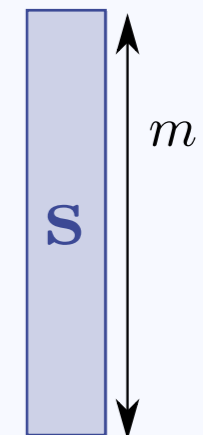
Subsampling



Coresets

Importance sampling

Here: linear « sketch »



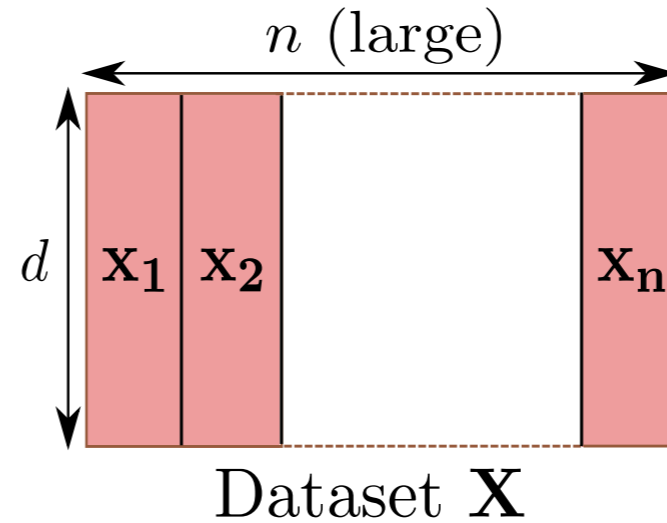
Only **one vector**

[Rémi Gribonval, Gilles Blanchard, Nicolas Keriven, Yann Traonmilin, Antoine Chatalic, Vincent Schellekens, Laurent Jacques...]

Theory of sketching

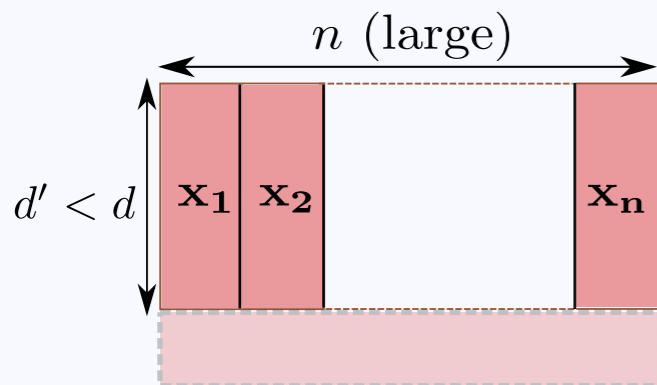
« Dimension » reduction

« Low-dim » representation of a dataset



How do we sketch ? How do we learn from sketch ?

Dimension reduction

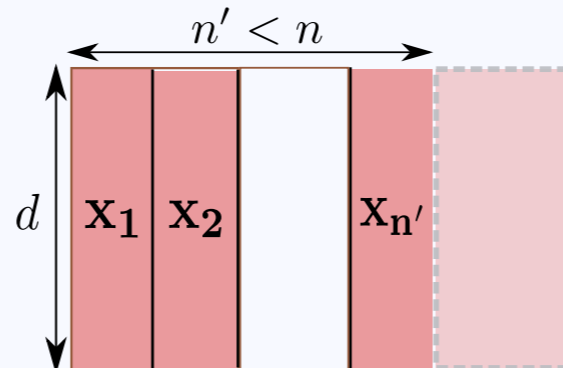


Random projections (JL lemma)

Feature selection

Minimum distortion embedding, PCA

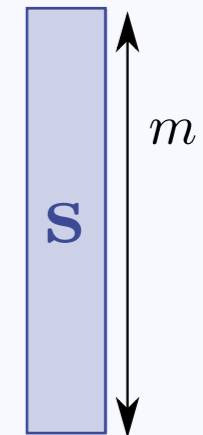
Subsampling



Coresets

Importance sampling

Here: linear « sketch »



Only one vector

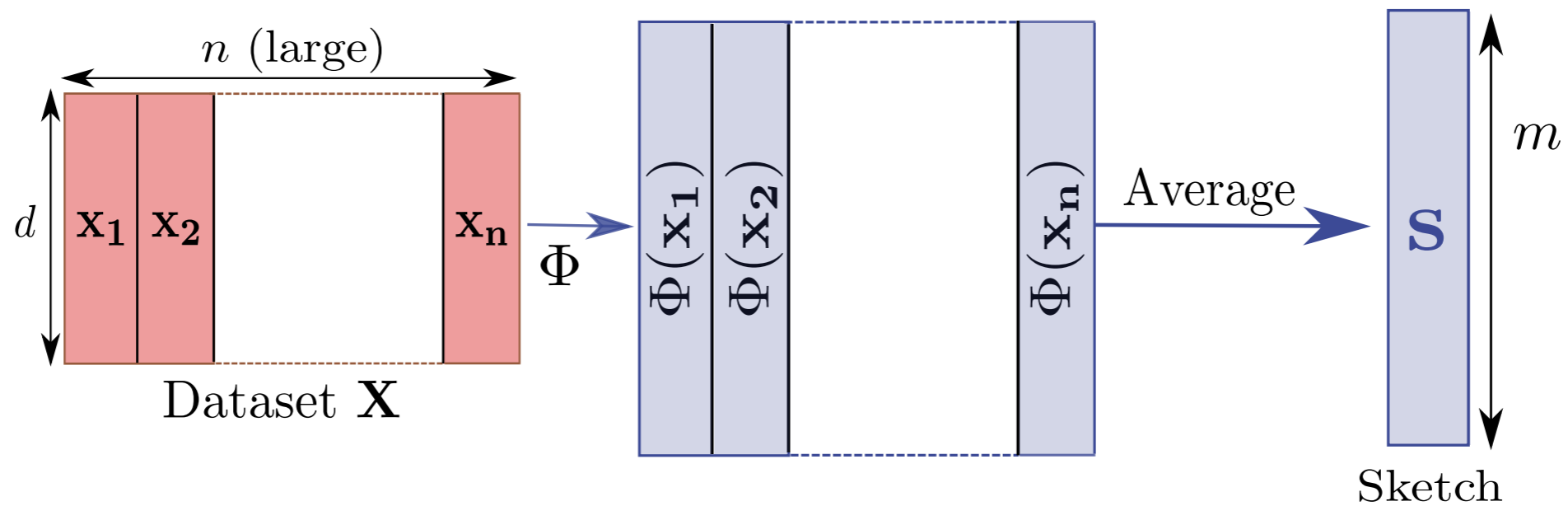
[Rémi Gribonval, Gilles Blanchard, Nicolas Keriven, Yann Traonmilin, Antoine Chatalic, Vincent Schellekens, Laurent Jacques...]

Theory of sketching

Obtaining the sketch

■ A function called **feature operator** $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

■ Averaging **n points** $\rightarrow \mathbf{S} := \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$

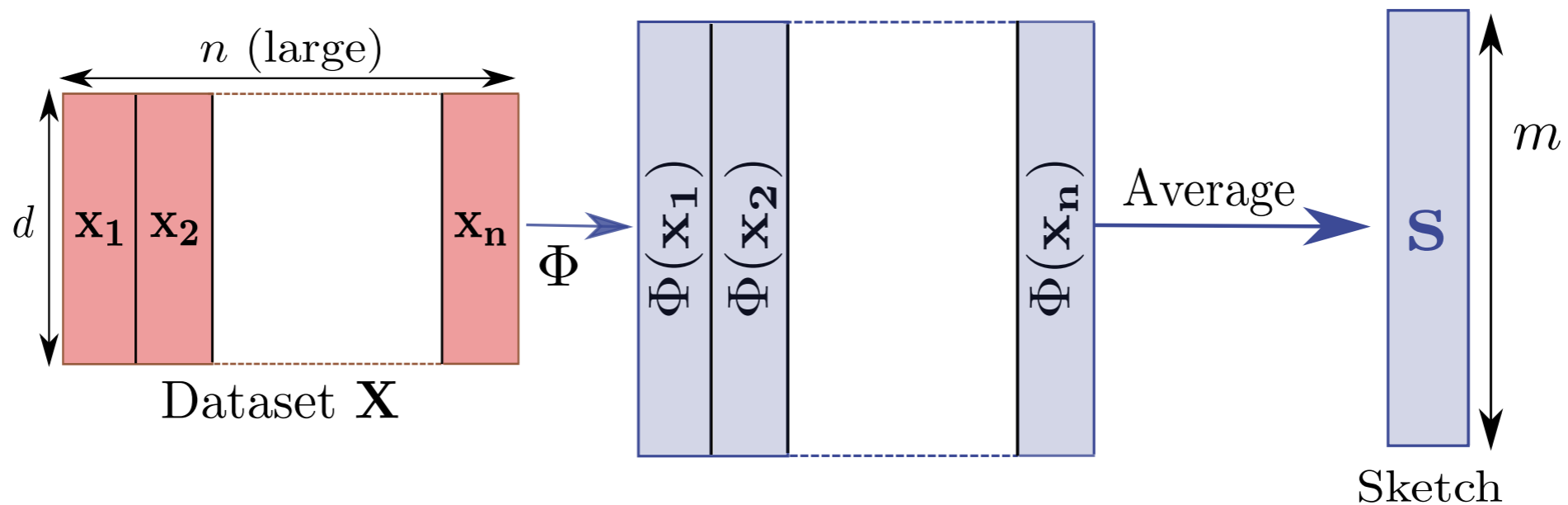


Theory of sketching

Obtaining the sketch

A function called **feature operator** $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

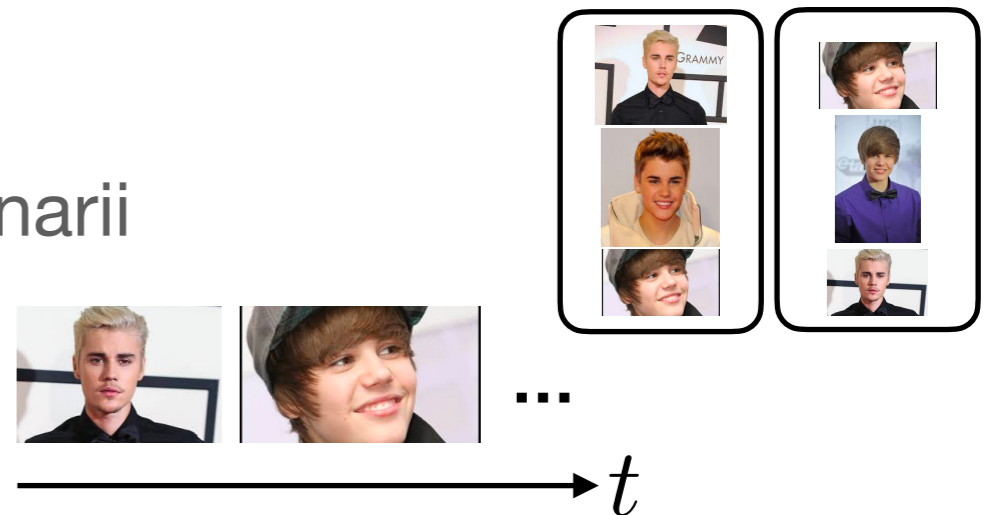
Averaging **n points** $\rightarrow \mathbf{s} := \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$



Average is a simple idea but ...

Suitable for **distributed / streaming** scenarios

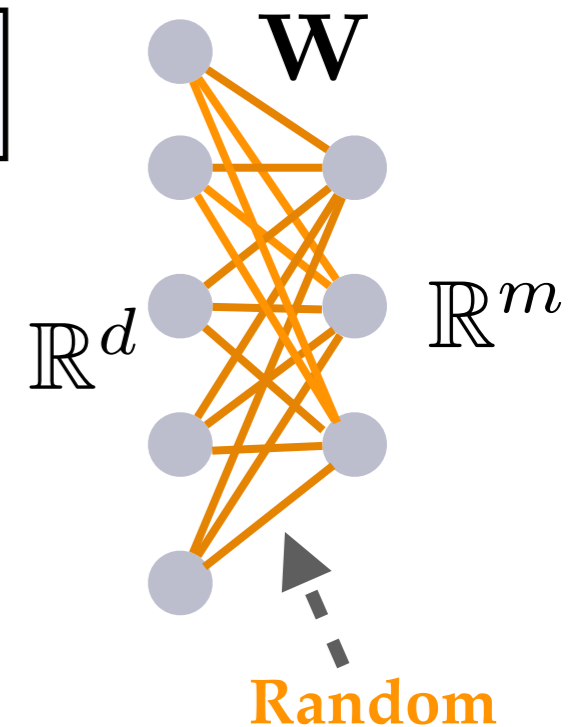
It can be calculated in **parallel**



| Theory of sketching

■ Randomization: the core of sketching

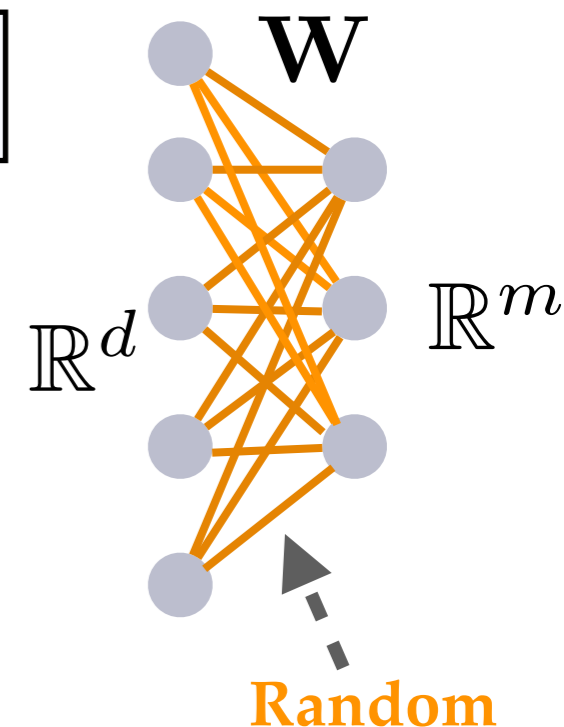
- A function called **feature operator** $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$
- In practice $\Phi(\mathbf{x}) = \rho(\mathbf{W}\mathbf{x})$
- $\mathbf{W} \in \mathbb{R}^{m \times d}$ is a **random matrix**
- $\rho : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is a **non-linear activation**



Theory of sketching

Randomization: the core of sketching

- A function called **feature operator** $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$
- In practice $\Phi(\mathbf{x}) = \rho(\mathbf{W}\mathbf{x})$
- $\mathbf{W} \in \mathbb{R}^{m \times d}$ is a **random matrix**
- $\rho : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is a **non-linear activation**



Random Fourier Features (RFF)

- $\mathbf{W} \in \mathbb{R}^{m \times d}$ is **Gaussian** $W_{ij} \sim \mathcal{N}(0, \sigma^2)$
- $\rho(\mathbf{y}) = \frac{1}{\sqrt{m}} (\exp(-iy_1), \dots, \exp(-iy_m))$

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{m}} (\exp(-i\omega_1^\top \mathbf{x}), \dots, \exp(-i\omega_m^\top \mathbf{x}))$$

$$\mathbf{W} = [\omega_1^\top, \dots, \omega_m^\top]$$

[Rahimi & Recht, 2008]



| Theory of sketching

■ A linear operator on distributions

■ The **sketching operator**

$$A : \pi \rightarrow \int_{\mathbb{R}^d} \Phi(\mathbf{x}) d\pi(\mathbf{x}) \in \mathbb{R}^m$$

■ This is a **linear operator** on measures/distributions

| Theory of sketching

■ A linear operator on distributions

- The **sketching operator**

$$\mathcal{A} : \pi \rightarrow \int_{\mathbb{R}^d} \Phi(\mathbf{x}) d\pi(\mathbf{x}) \in \mathbb{R}^m$$

- This is a **linear operator** on measures/distributions

- **Generalized moments** $\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\Phi(\mathbf{x})]$

Theory of sketching

A linear operator on distributions

- The **sketching operator**

$$\mathcal{A} : \pi \rightarrow \int_{\mathbb{R}^d} \Phi(\mathbf{x}) d\pi(\mathbf{x}) \in \mathbb{R}^m$$

- This is a **linear operator** on measures/distributions

Mean (moment 1)

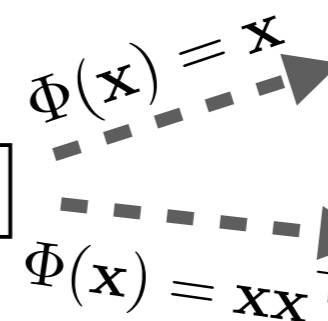
$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\mathbf{x}]$$

- Generalized moments

$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\Phi(\mathbf{x})]$$

Variance (moment 2)

$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\mathbf{x}\mathbf{x}^\top]$$



Theory of sketching

A linear operator on distributions

- The **sketching operator**

$$\mathcal{A} : \pi \rightarrow \int_{\mathbb{R}^d} \Phi(\mathbf{x}) d\pi(\mathbf{x}) \in \mathbb{R}^m$$

- This is a **linear operator** on measures/distributions

Mean (moment 1)

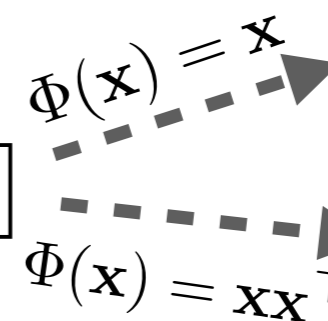
$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\mathbf{x}]$$

- Generalized moments

$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\Phi(\mathbf{x})]$$

Variance (moment 2)

$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\mathbf{x}\mathbf{x}^\top]$$



$$\mathcal{A}\pi = (\mathcal{F}[\pi](\omega_1), \dots, \mathcal{F}[\pi](\omega_m))$$

Sampling of the **Fourier transform** of the distrib.

Theory of sketching

A linear operator on distributions

- The **sketching operator**

$$\mathcal{A} : \pi \rightarrow \int_{\mathbb{R}^d} \Phi(\mathbf{x}) d\pi(\mathbf{x}) \in \mathbb{R}^m$$

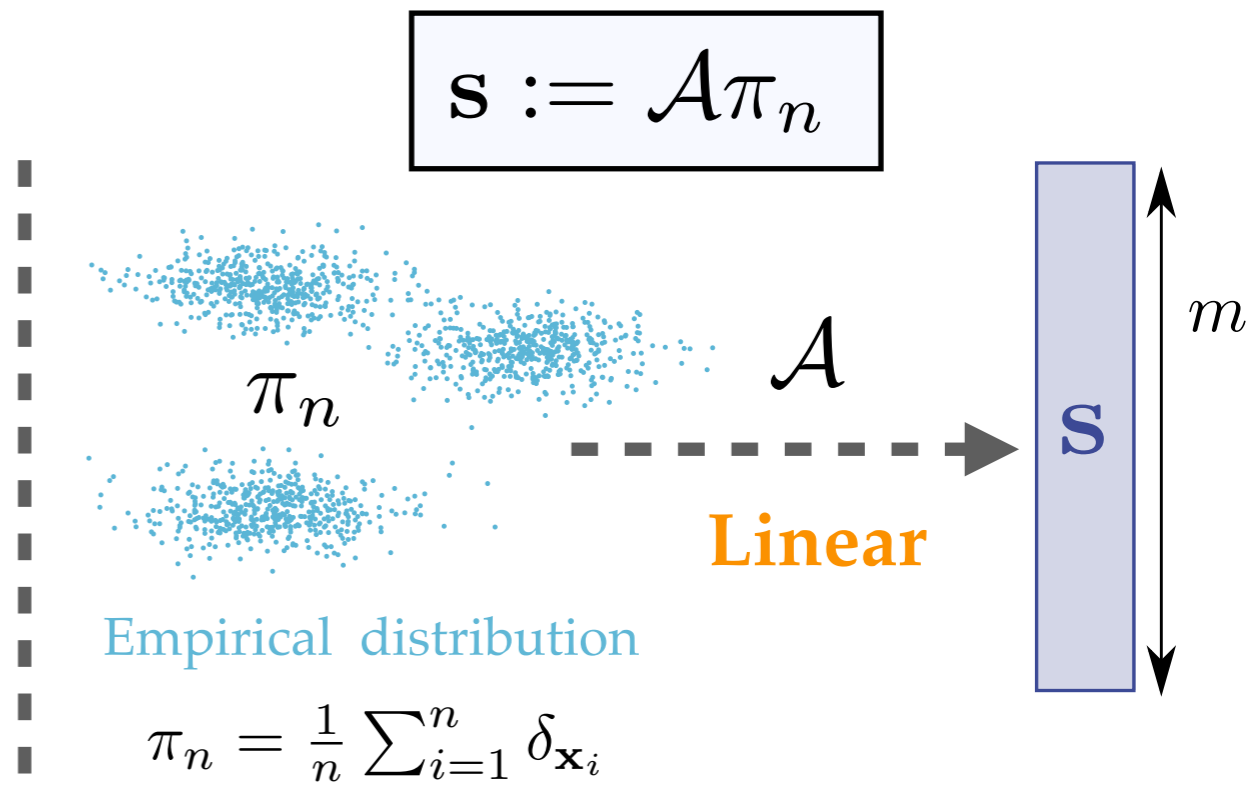
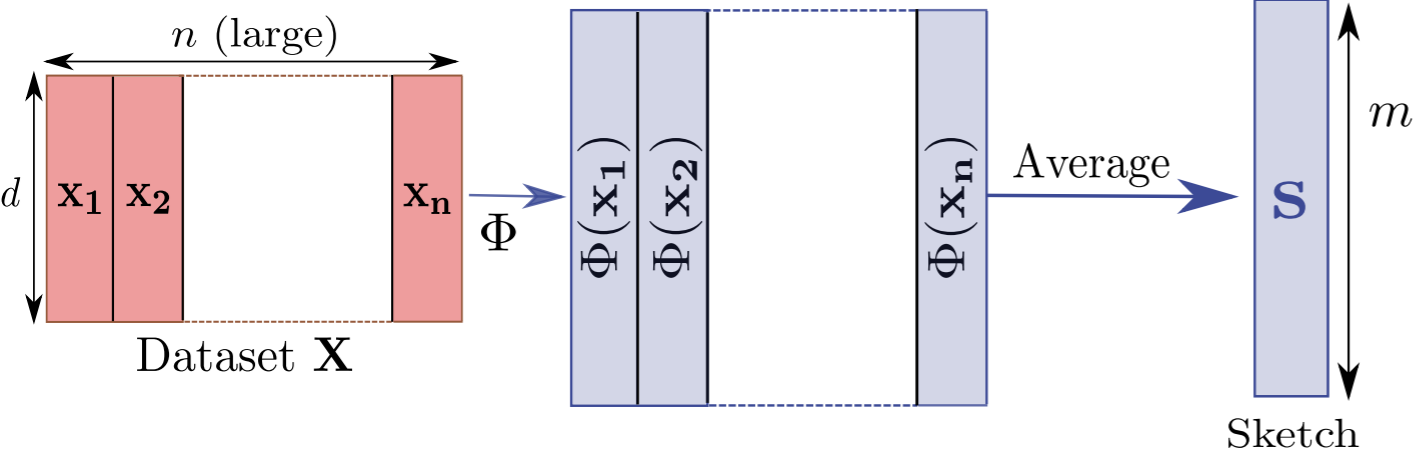
- This is a **linear operator** on measures/distributions

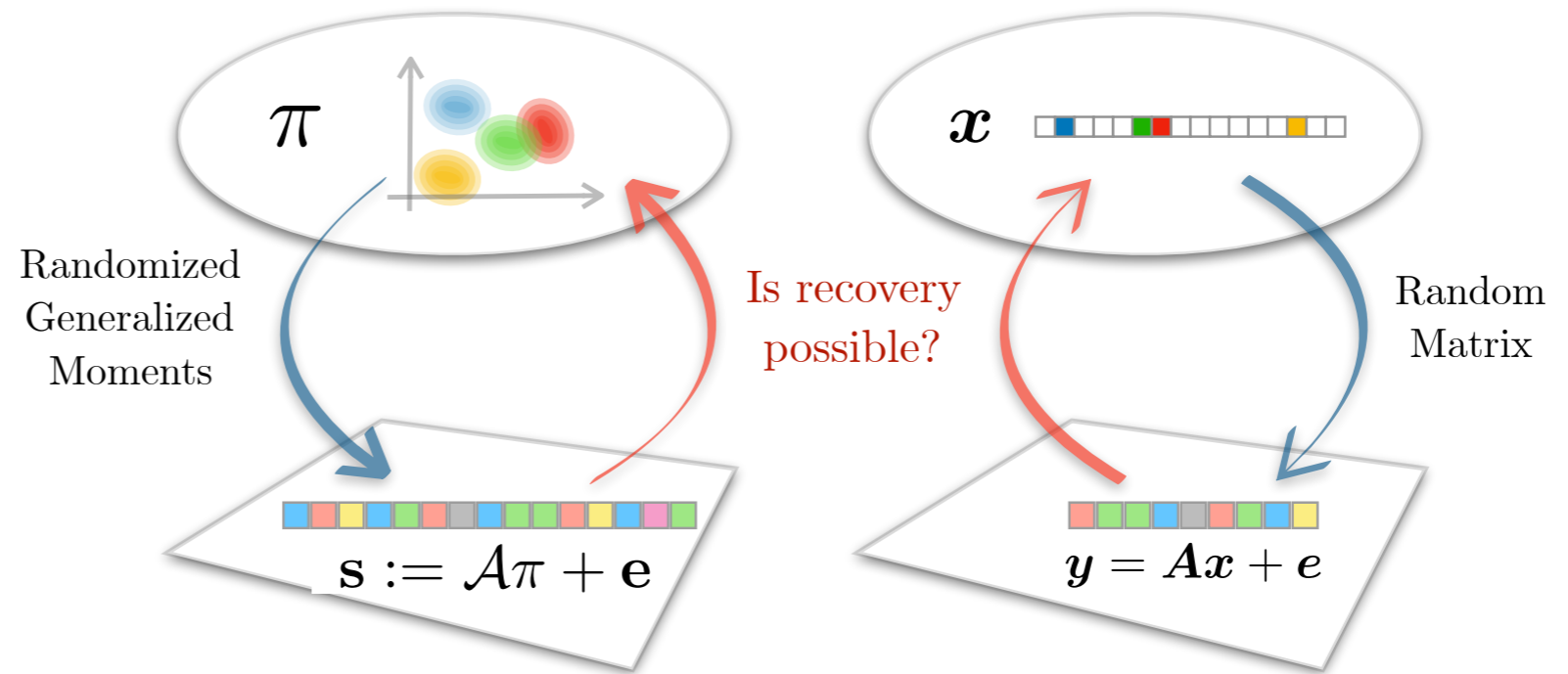
- Generalized moments $\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\Phi(\mathbf{x})]$

The two « ways » of sketching

$$\mathbf{s} := \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$$

$$\mathbf{s} := \mathcal{A}\pi_n$$





Compressive Learning

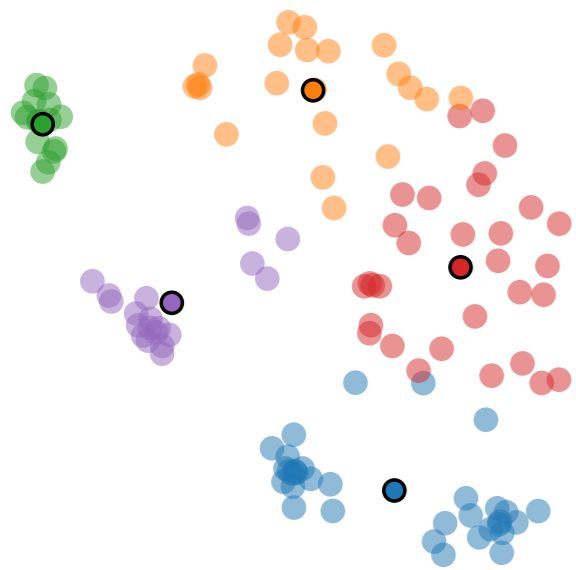
- Theory of sketching
- Sketching in practice
- RIP for theoretical guarantees
- Limitations & perspectives

Sketching in practice

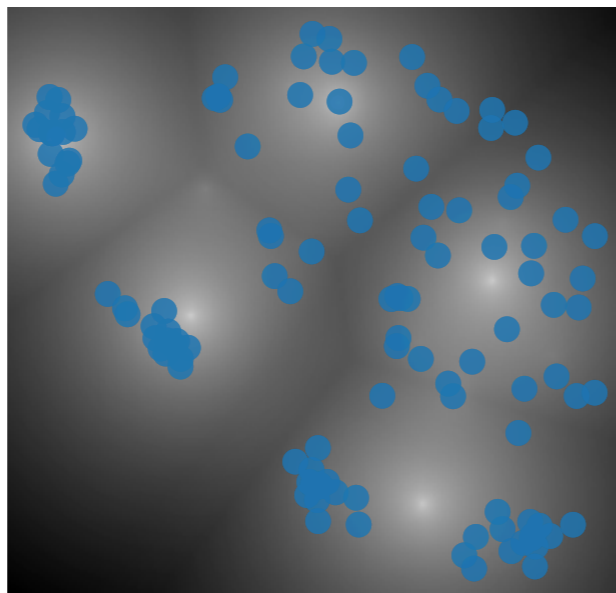
■ Come back to K-means:

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \frac{1}{n} \sum_{i=1}^n \min_{k \in \llbracket K \rrbracket} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

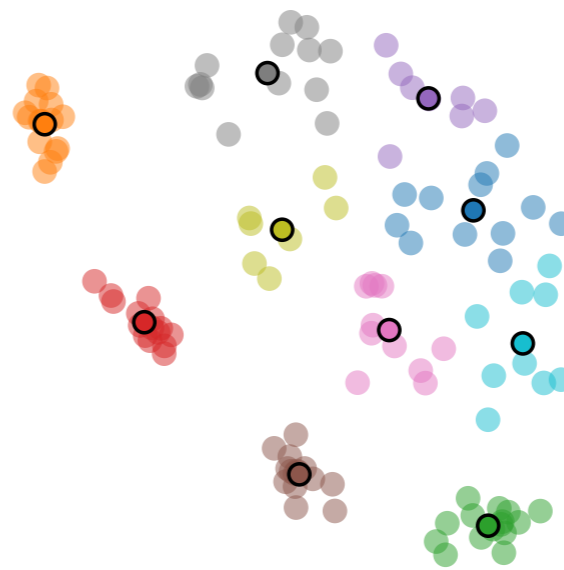
K-means for K=5



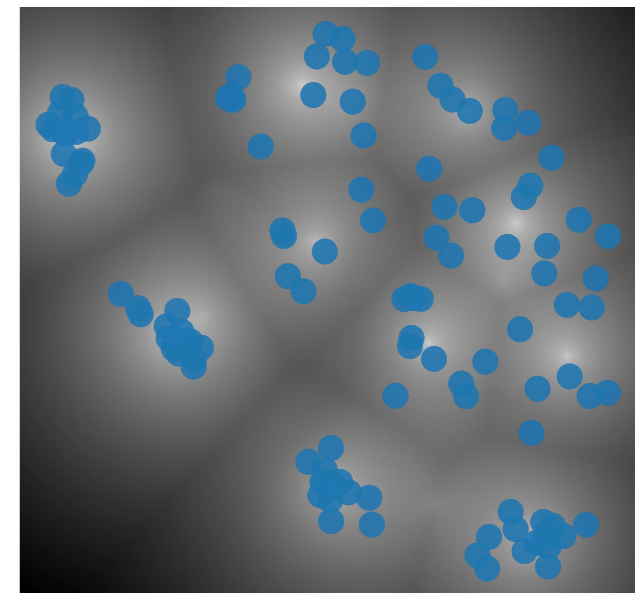
$f(x) = \min_k |x - c_k|^2$ for K=5



K-means for K=10



$f(x) = \min_k |x - c_k|^2$ for K=10



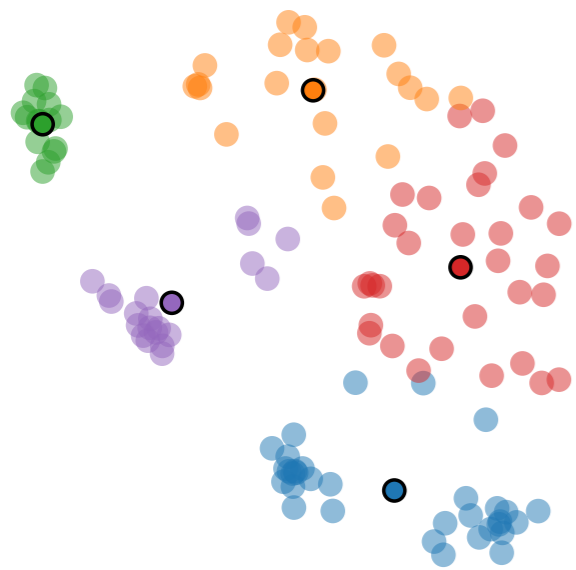
Sketching in practice

We want to find Kd params.

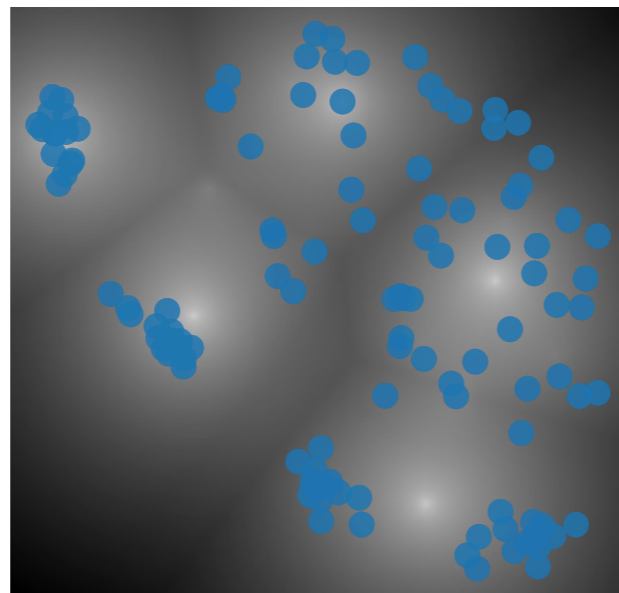
■ Come back to K-means:

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \frac{1}{n} \sum_{i=1}^n \min_{k \in \llbracket K \rrbracket} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

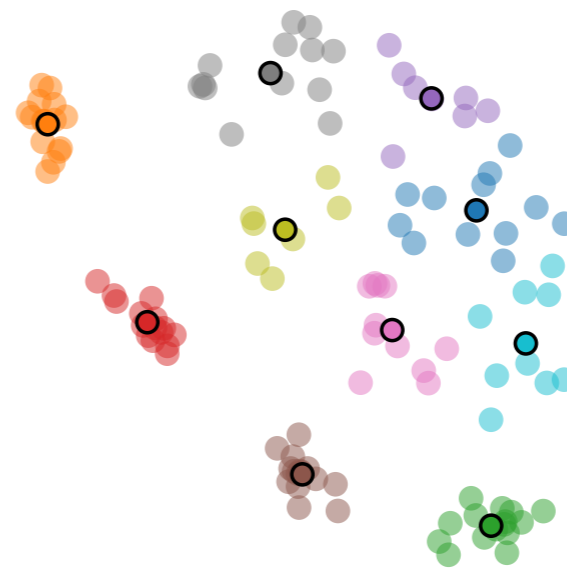
K-means for K=5



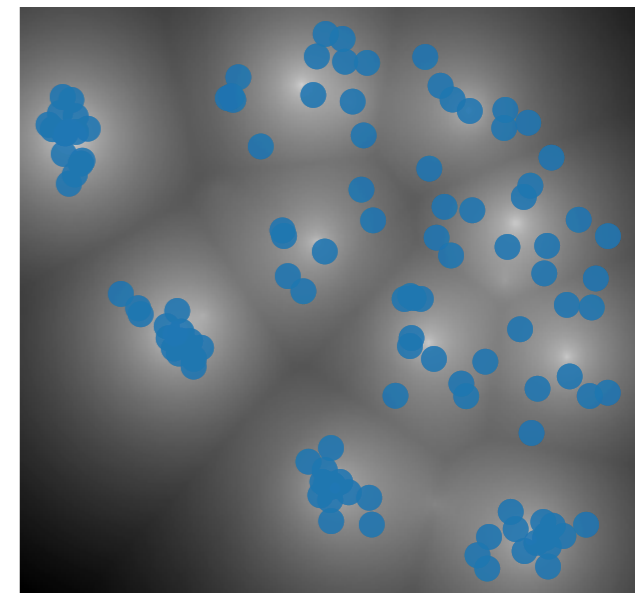
$f(x) = \min_k |x - c_k|^2$ for K=5



K-means for K=10



$f(x) = \min_k |x - c_k|^2$ for K=10

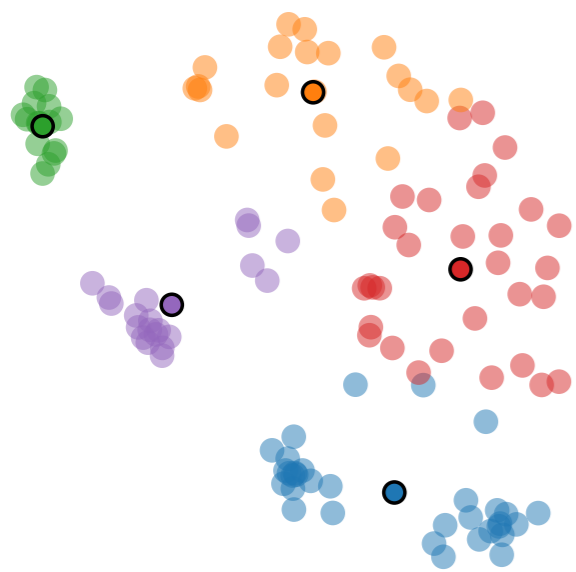


Sketching in practice

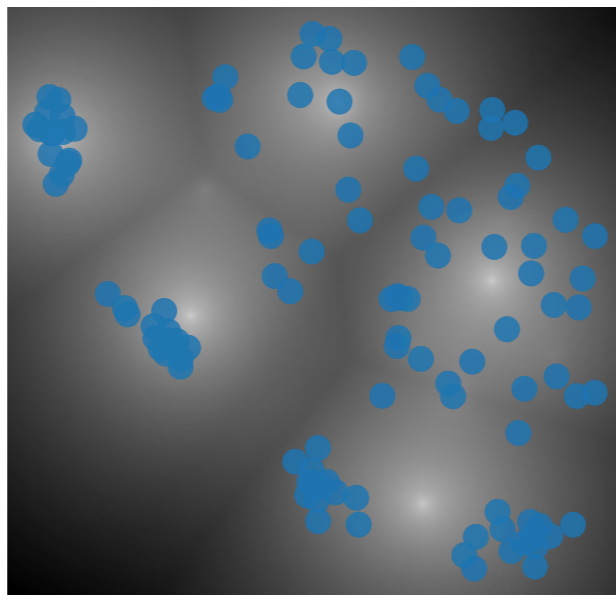
■ Come back to K-means:

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \frac{1}{n} \sum_{i=1}^n \min_{k \in \llbracket K \rrbracket} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

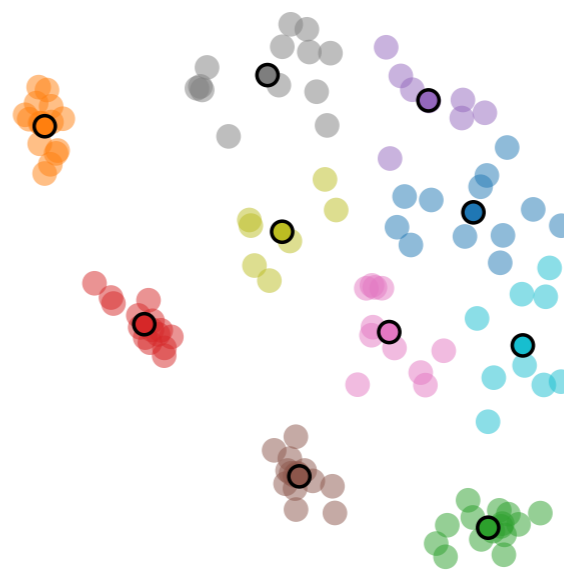
K-means for K=5



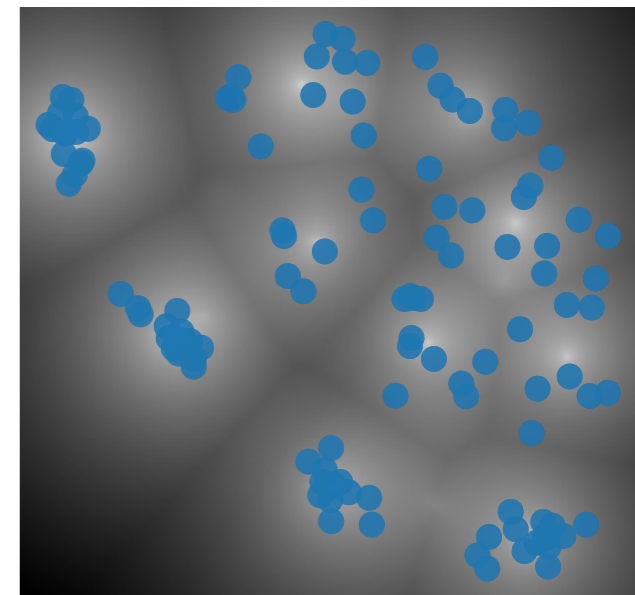
$f(x) = \min_k |x - c_k|^2$ for K=5



K-means for K=10



$f(x) = \min_k |x - c_k|^2$ for K=10



■ Another point of view:

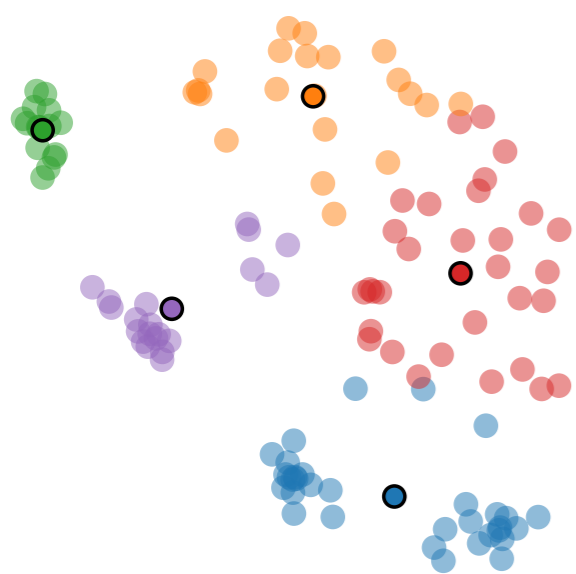
- Find a **distribution** $\hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k}$ that **bests approximate** π
- We have access to the whole dataset i.e. π_n

Sketching in practice

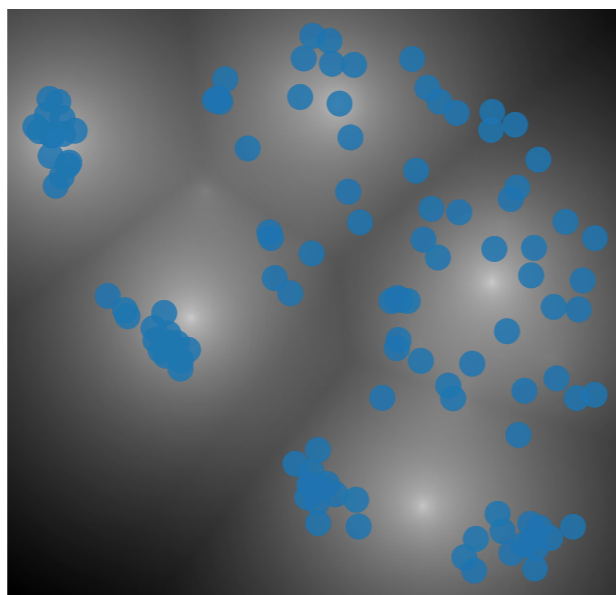
Come back to K-means:

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \frac{1}{n} \sum_{i=1}^n \min_{k \in \llbracket K \rrbracket} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

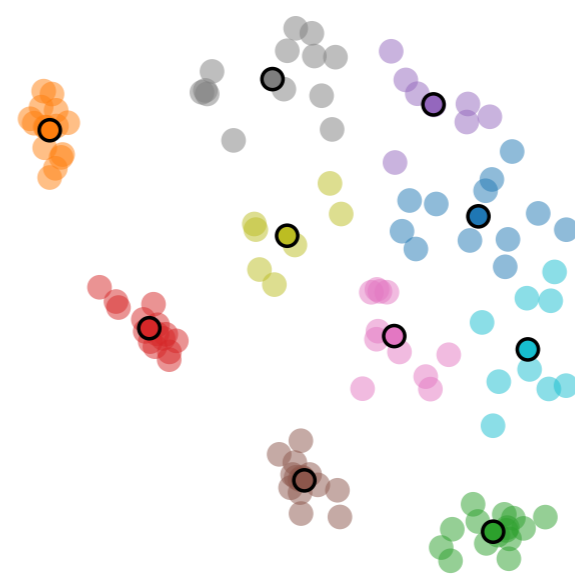
K-means for K=5



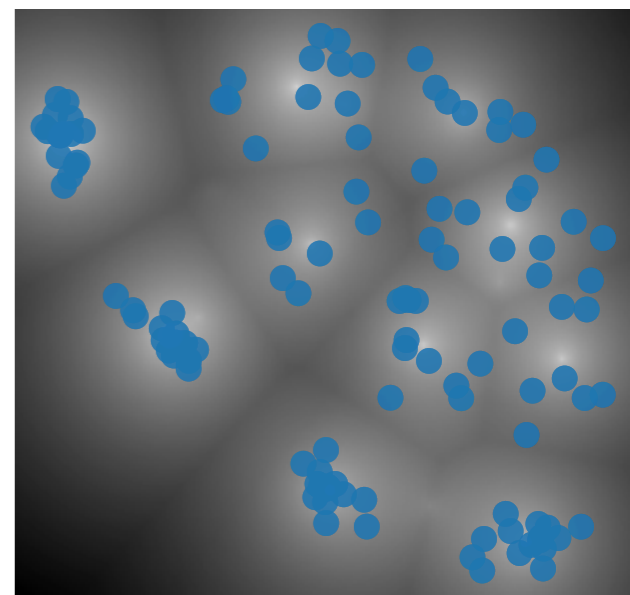
$f(x) = \min_k |x - c_k|^2$ for K=5



K-means for K=10



$f(x) = \min_k |x - c_k|^2$ for K=10



Another point of view:

- Find a **distribution** $\hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k}$ that **bests approximate** π
- We have access to the whole dataset i.e. π_n

In sketching:

- We only have access to

$$\mathbf{s} := \mathcal{A}\pi_n$$

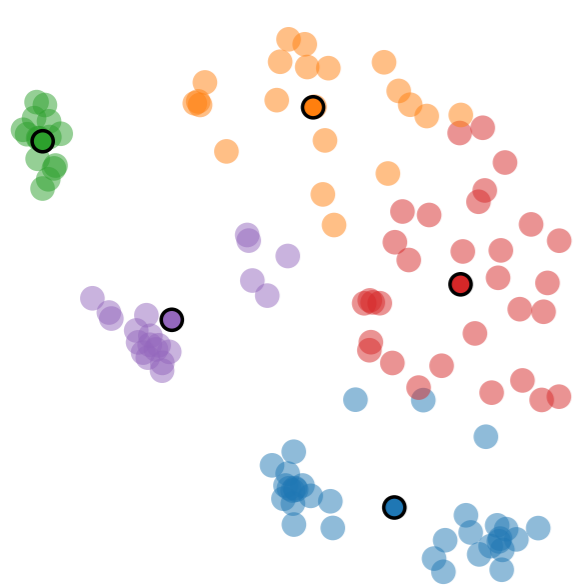
How do we do ?

Sketching in practice

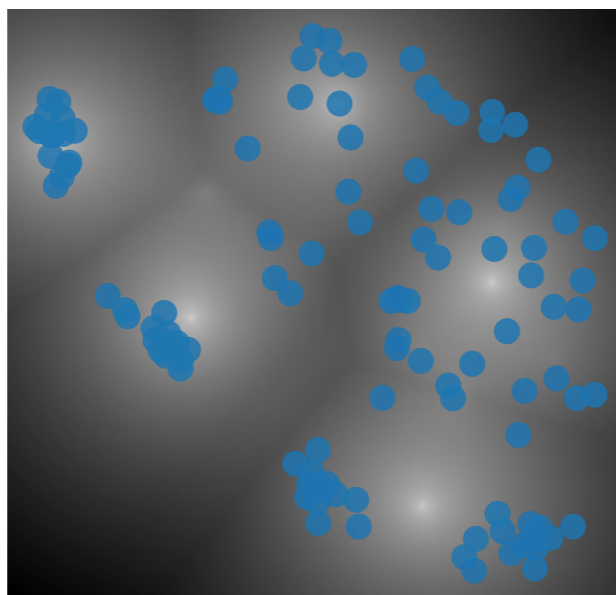
Come back to K-means:

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \frac{1}{n} \sum_{i=1}^n \min_{k \in \llbracket K \rrbracket} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

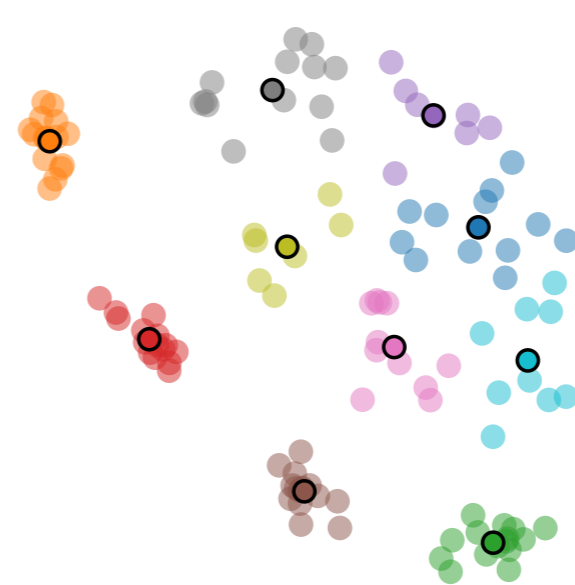
K-means for K=5



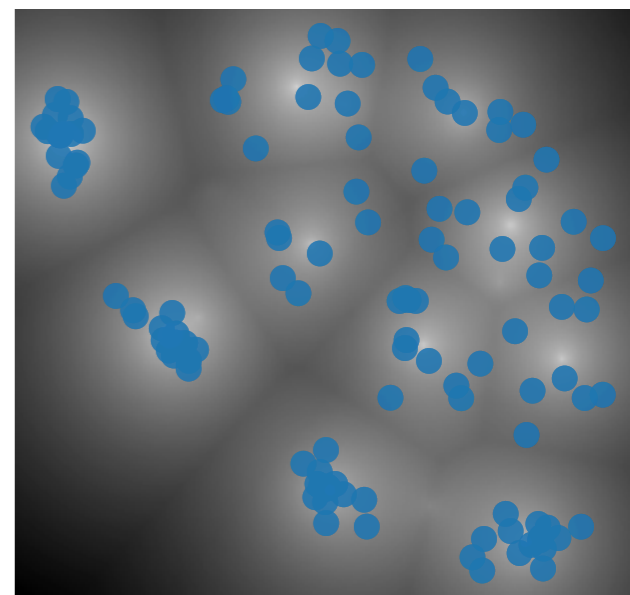
$f(x) = \min_k |x - c_k|^2$ for K=5



K-means for K=10



$f(x) = \min_k |x - c_k|^2$ for K=10



Another point of view:

Find a **distribution** $\hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k}$ that **bests approximate** π

We have access to the whole dataset i.e. π_n

we will change that

In sketching:

We only have access to

$$\mathbf{s} := \mathcal{A}\pi_n$$

How do we do ?

Sketching in practice

■ Sketching for large scale K-means:

■ We aim at solving:

$$\hat{\pi} \text{ s.t. } \hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k} \quad \min_{\hat{\pi}} \|\mathbf{s} - \mathcal{A}\hat{\pi}\|_2$$

Sketching in practice

■ Sketching for large scale K-means:

■ We aim at solving:

$$\min_{\hat{\pi}} \|\mathbf{s} - \mathcal{A}\hat{\pi}\|_2$$

$\hat{\pi}$ s.t. $\hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k}$

Find a distrib. of K diracs

Sketching in practice

■ Sketching for large scale K-means:

■ We aim at solving:

$$\hat{\pi} \text{ s.t. } \hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k} \quad \min \|\mathbf{s} - \mathcal{A}\hat{\pi}\|_2$$

Sketch of the distrib

Sketching in practice

■ Sketching for large scale K-means:

■ We aim at solving:

$$\hat{\pi} \text{ s.t. } \hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k} \quad \min \|\mathbf{s} - \mathcal{A}\hat{\pi}\|_2$$

Sketch of the data

Sketching in practice

■ Sketching for large scale K-means:

- We aim at solving:

$$\hat{\pi} \text{ s.t. } \hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k} \quad \min_{\hat{\pi}} \|\mathbf{s} - \mathcal{A}\hat{\pi}\|_2$$

- Find a distribution of K diracs whose sketch is the **closest** to the **sketch of the dataset**
- **Different criteria** than K-means

Sketching in practice

Sketching for large scale K-means:

- We aim at solving:

$$\hat{\pi} \text{ s.t. } \hat{\pi} = \frac{1}{K} \sum_{k=1}^K \delta_{\mathbf{c}_k} \quad \min_{\hat{\pi}} \|\mathbf{s} - \mathcal{A}\hat{\pi}\|_2$$

- Find a distribution of K diracs whose sketch is the closest to the sketch of the dataset
- Different criteria than K-means

Reformulation:

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \left\| \mathbf{s} - \frac{1}{K} \sum_{k=1}^K \Phi(\mathbf{c}_k) \right\|_2$$

Sketching in practice

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \left\| \mathbf{s} - \frac{1}{K} \sum_{k=1}^K \Phi(\mathbf{c}_k) \right\|_2$$

Algorithm:

- Inspired from orthogonal matching pursuit (OMP)

Sketching in practice

$$\min_{\mathbf{c}_1, \dots, \mathbf{c}_K} \left\| \mathbf{s} - \frac{1}{K} \sum_{k=1}^K \Phi(\mathbf{c}_k) \right\|_2$$

Algorithm:

Inspired from orthogonal matching pursuit (OMP)

$\Theta \leftarrow \emptyset$ $\mathbf{r} \leftarrow \mathbf{s}$ // Initialize

while $|\Theta| \leq K$:

$\hat{\mathbf{c}} \in \arg \max_{\mathbf{c} \in \mathbb{R}^d} \left| \left\langle \frac{\Phi(\mathbf{c})}{\|\Phi(\mathbf{c})\|}, \mathbf{r} \right\rangle \right|$ // Find a new atom:
minimizes the residuals

$\Theta \leftarrow \Theta \cup \{\hat{\mathbf{c}}\}$ // add it to the support $\Theta = (\mathbf{c}_1, \dots, \mathbf{c}_{|\Theta|})$

$\hat{\alpha} \in \arg \min_{\alpha_1, \dots, \alpha_k} \left\| \mathbf{s} - \sum_{k=1}^{|\Theta|} \alpha_k \Phi(\mathbf{c}_k) \right\|^2$ // Adjust weights:
least-squares

$\mathbf{r} \leftarrow \mathbf{s} - \sum_{k=1}^{|\Theta|} \hat{\alpha}_k \Phi(\mathbf{c}_k)$ // update residuals

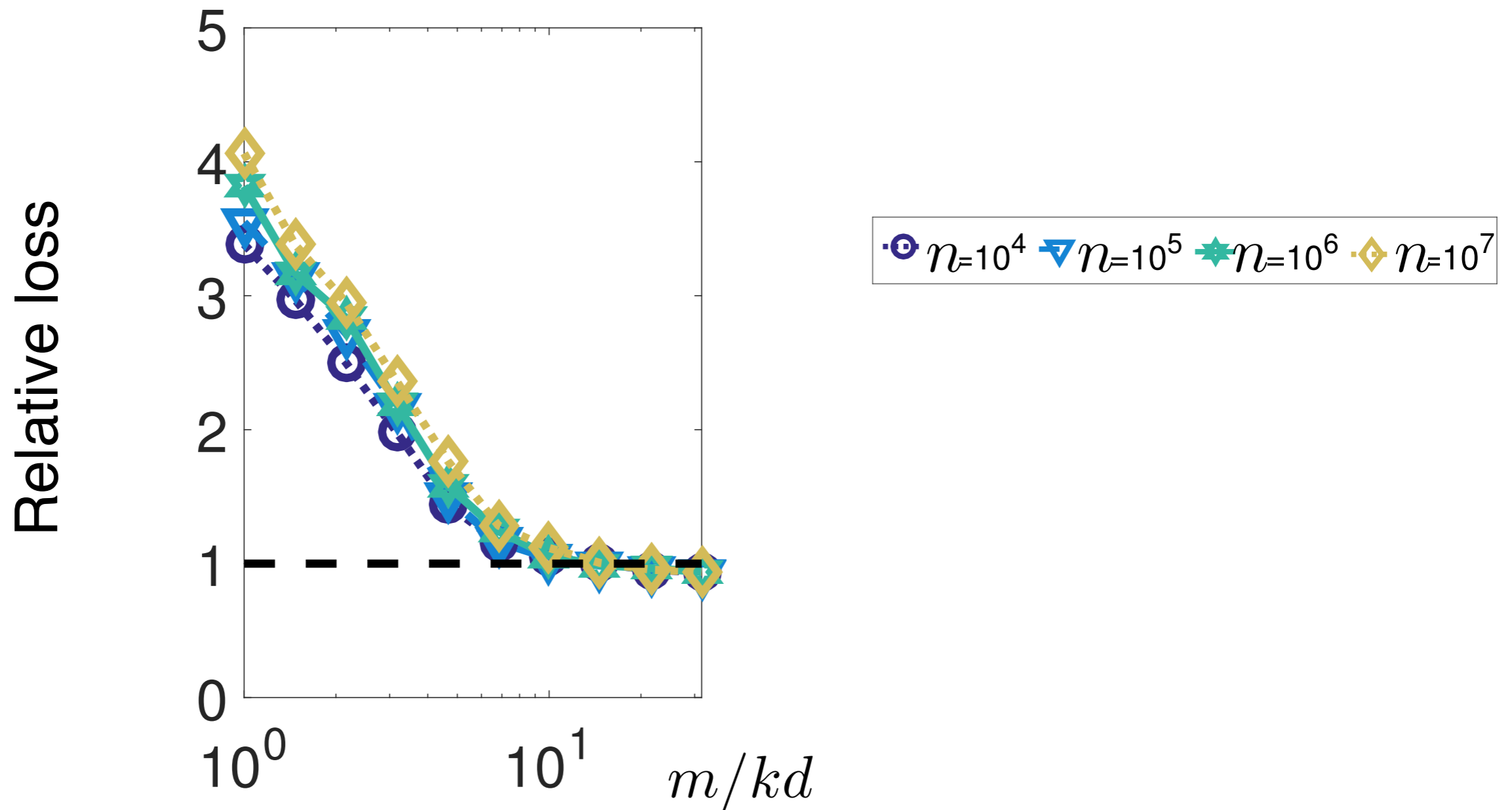
return: $(\mathbf{c}_1, \dots, \mathbf{c}_{|\Theta|})$

Complexity: $\mathcal{O}(mdK^2)$

Sketching in practice

■ **Results: How do we choose m ?** $m \approx Kd$ **number of params?**

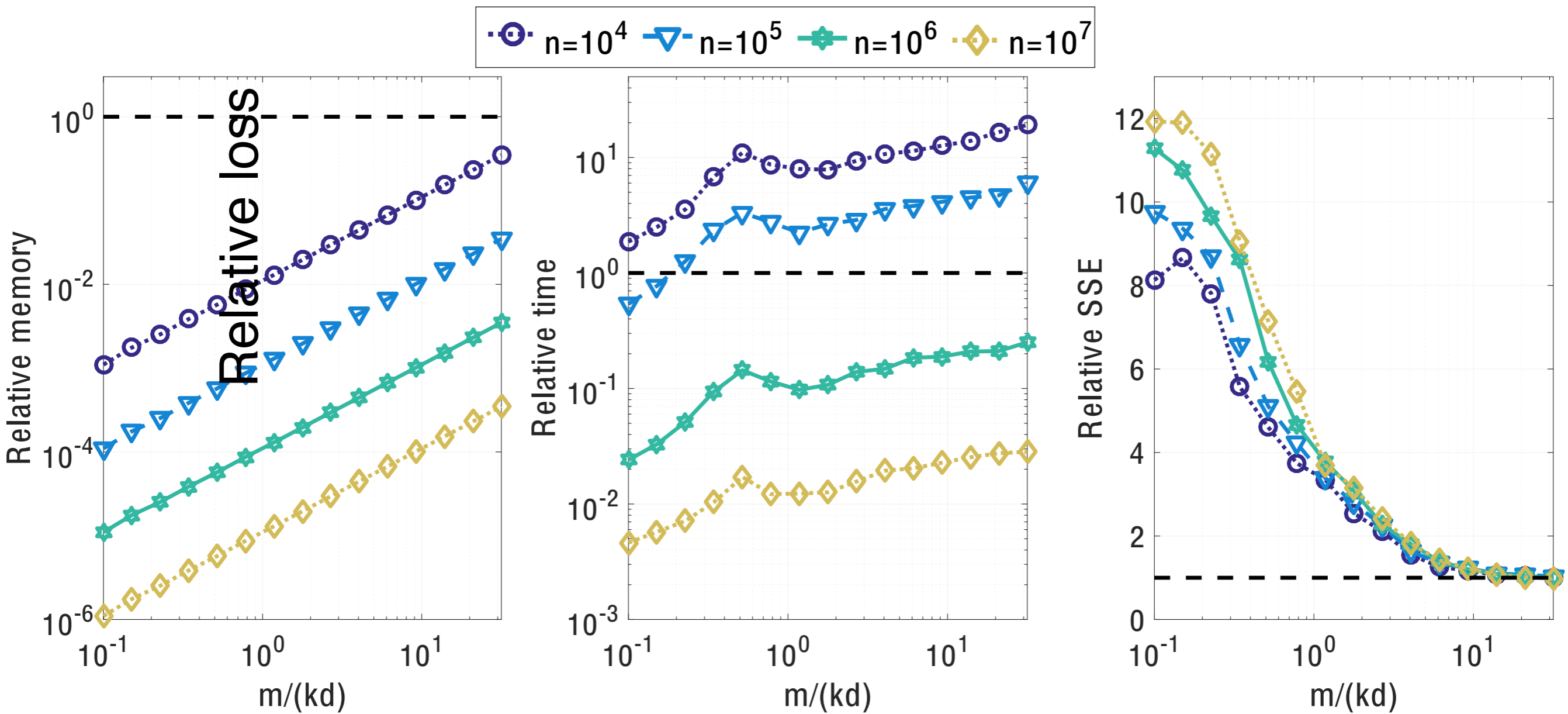
■ Synthetic dataset: $d = 10, K = 10$



Sketching in practice

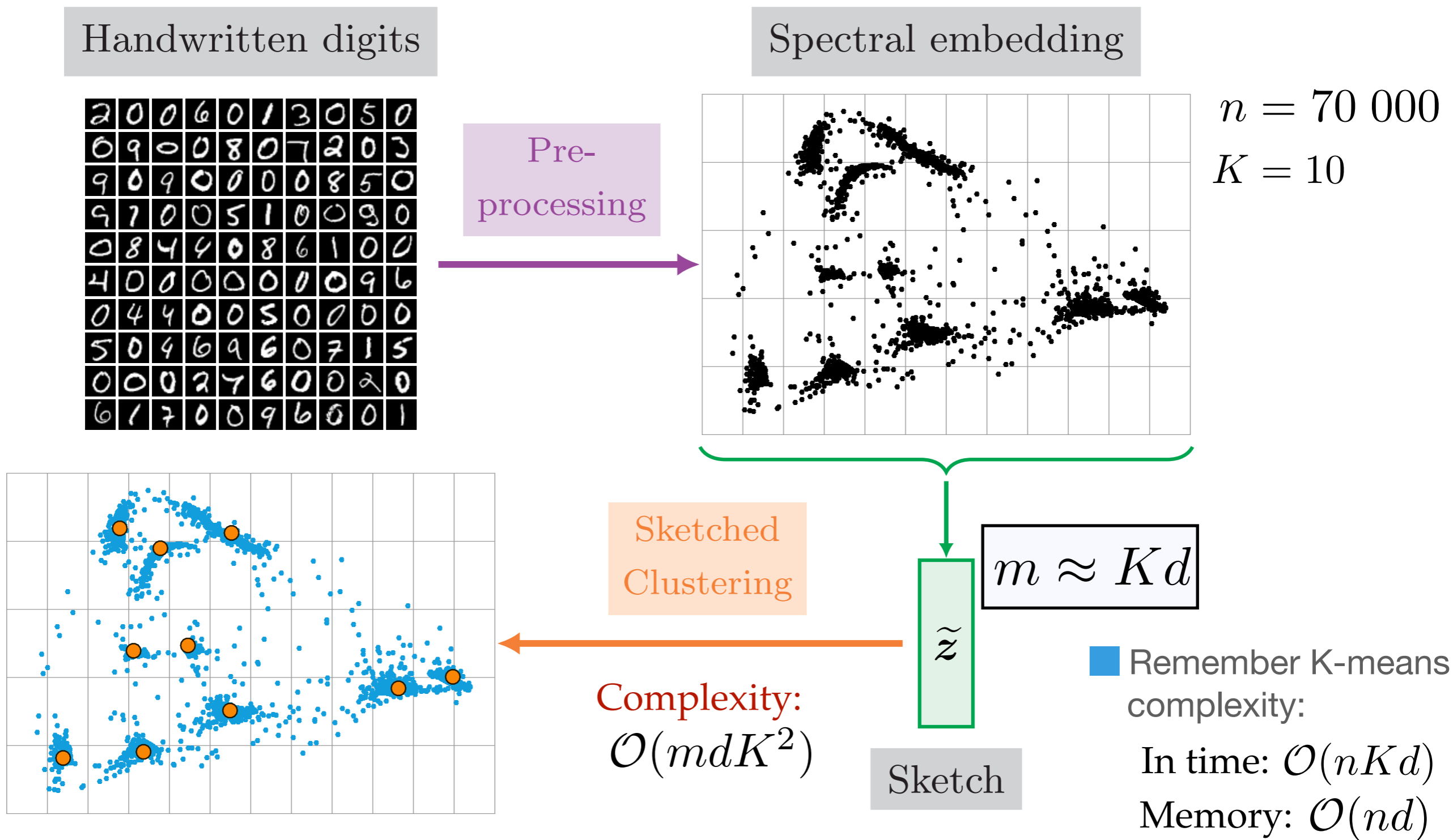
Results: How do we choose m ? $m \approx Kd$ number of params?

■ Synthetic dataset: $d = 10, K = 10$



Sketching in practice

Results for K-means:



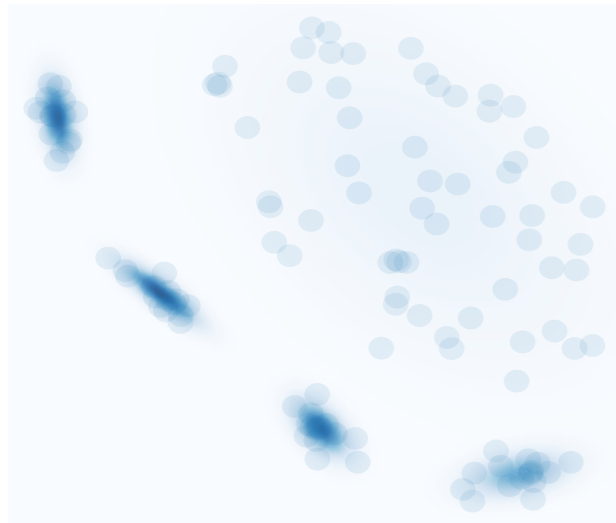
Sketching in practice

GMM:
$$\pi_{\theta}(\mathbf{x}) = \sum_{k=1}^K \alpha_k \pi_{\mu_k, \Sigma_k}(\mathbf{x})$$

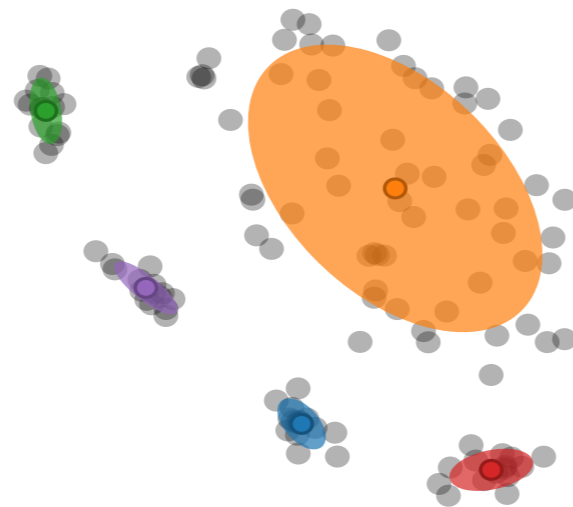
■ Come back to GMM:

MLE estimate:
$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n -\log(\pi_{\theta}(\mathbf{x}_i))$$

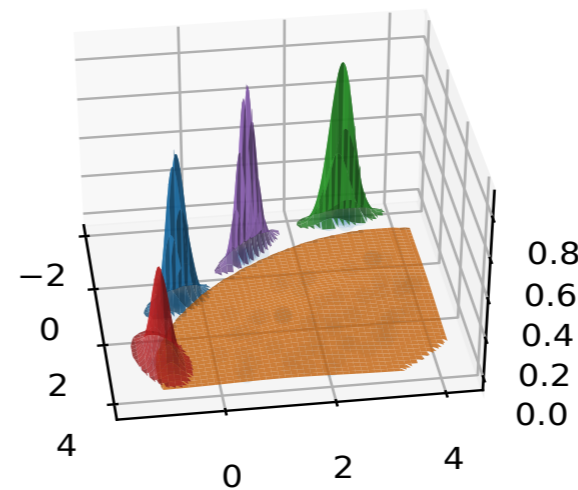
GMM density



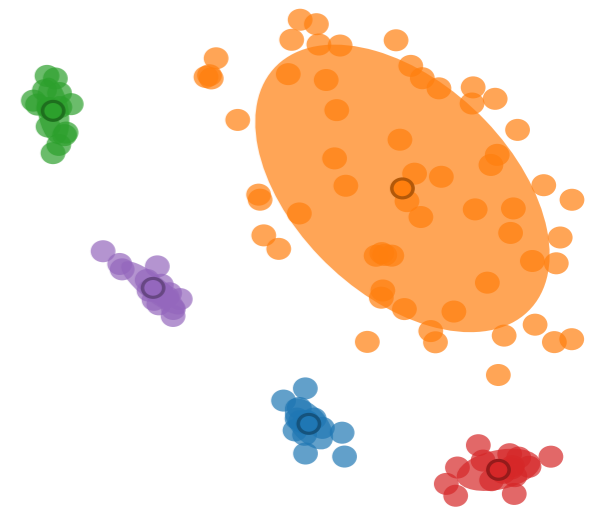
Estimated GMM



GMM mixture densities



GMM clustering



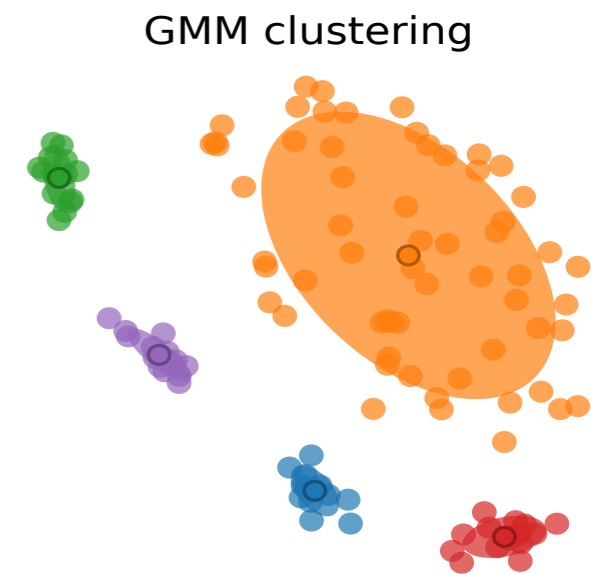
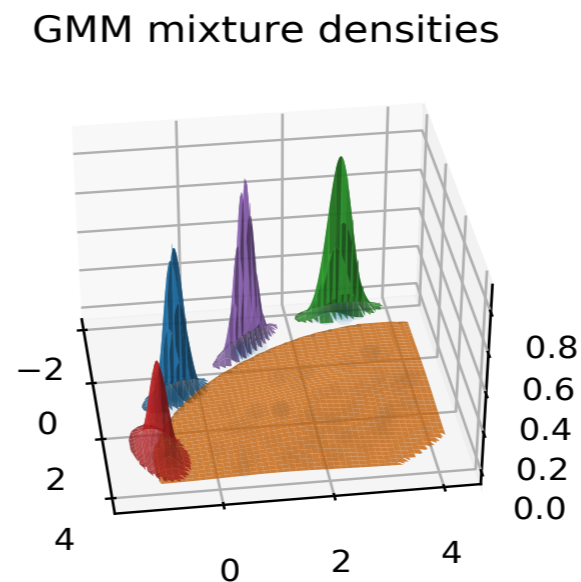
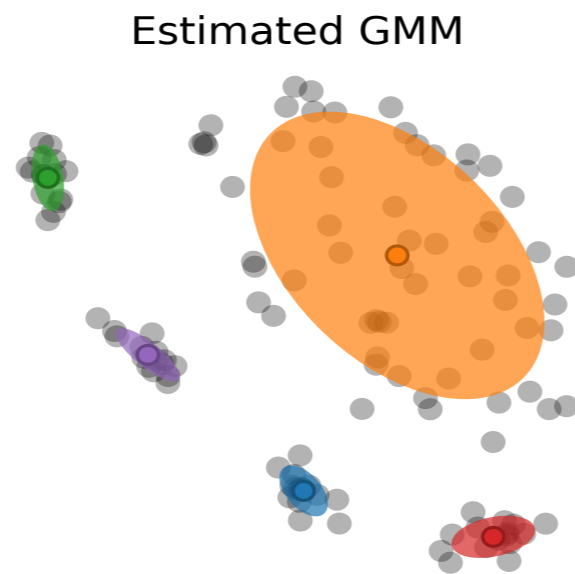
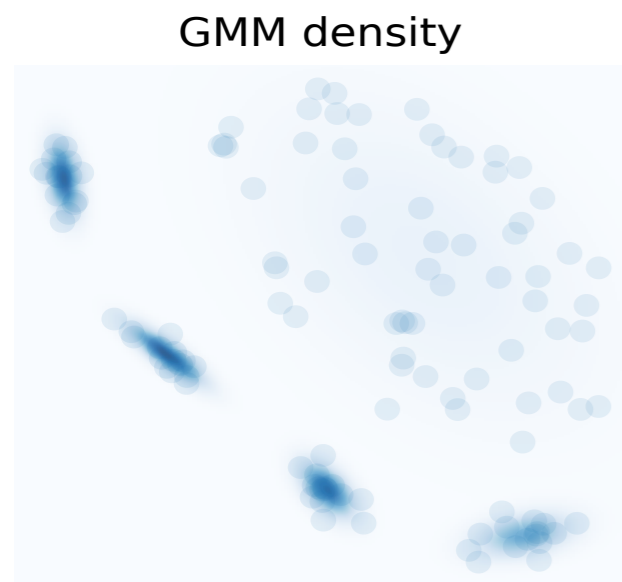
■ Find a **distribution** π_{θ} that **bests approximate** π

Sketching in practice

GMM:
$$\pi_{\theta}(\mathbf{x}) = \sum_{k=1}^K \alpha_k \pi_{\mu_k, \Sigma_k}(\mathbf{x})$$

■ Come back to GMM:

MLE estimate:
$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n -\log(\pi_{\theta}(\mathbf{x}_i))$$



■ Find a **distribution** π_{θ} that **bests approximate** π


■ Same idea than before:

$$\min_{\theta = \{\alpha_k, \mu_k, \Sigma_k\}_{k=1}^K} \|\mathbf{s} - \mathcal{A}\pi_{\theta}\|_2$$

Find a GMM **whose sketch** is the **closest** to the **sketch of the dataset**

Sketching in practice

■ Sketching for large scale GMM:

- A little bit more delicate: need to evaluate $\mathcal{A}\pi_{\theta}$
- Linearity: $\mathcal{A}\pi_{\theta} = \sum_{k=1}^K \alpha_k \mathcal{A}(\pi_{\mu_k, \Sigma_k})$ 
- When $\Phi = \text{RFF}$:

$\mathcal{A}(\pi_{\mu_k, \Sigma_k})$ is the **Fourier transform of a Gaussian** evaluated at m points

But **Fourier transform of a Gaussian = Gaussian**

-> Just need to sample m points from a Gaussian (**easy**)

Sketching in practice

Sketching for large scale GMM:

- A little bit more delicate: need to evaluate $\mathcal{A}\pi_{\theta}$
- Linearity: $\mathcal{A}\pi_{\theta} = \sum_{k=1}^K \alpha_k \mathcal{A}(\pi_{\mu_k, \Sigma_k})$
- When $\Phi = \text{RFF}$:

GMM

$\mathcal{A}(\pi_{\mu_k, \Sigma_k})$ is the **Fourier transform of a Gaussian** evaluated at m points

But **Fourier transform of a Gaussian = Gaussian**

-> Just need to sample m points from a Gaussian (**easy**)

Algorithm:

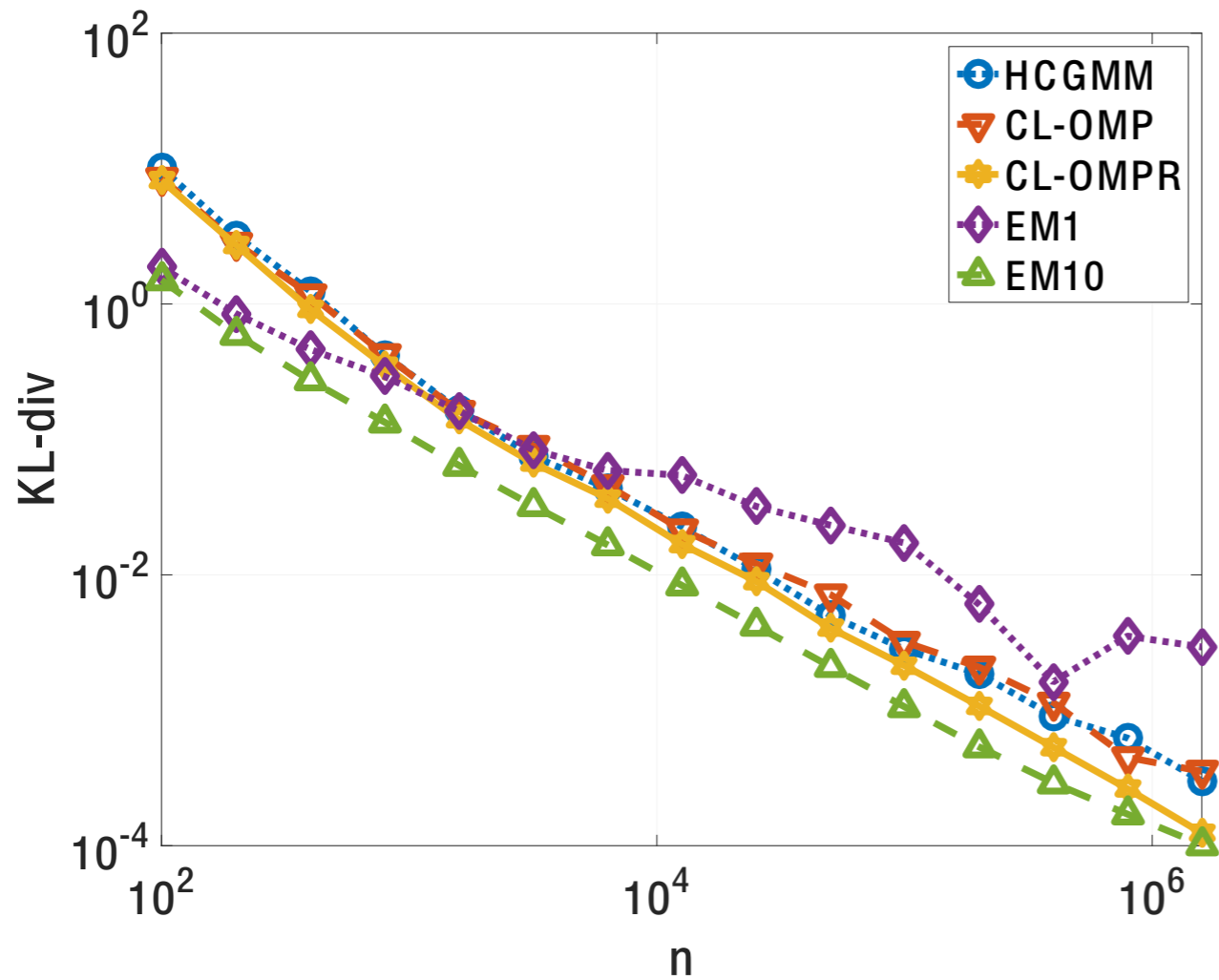
- The algorithm of K-means adapts to this setting

$$\hat{\mu}, \hat{\Sigma} \leftarrow \arg \max_{\mu, \Sigma} \left\langle \frac{\mathcal{A}\pi_{\mu, \Sigma}}{\|\mathcal{A}\pi_{\mu, \Sigma}\|}, \mathbf{r} \right\rangle \quad // \text{ Find a new atom: minimizes the residuals}$$

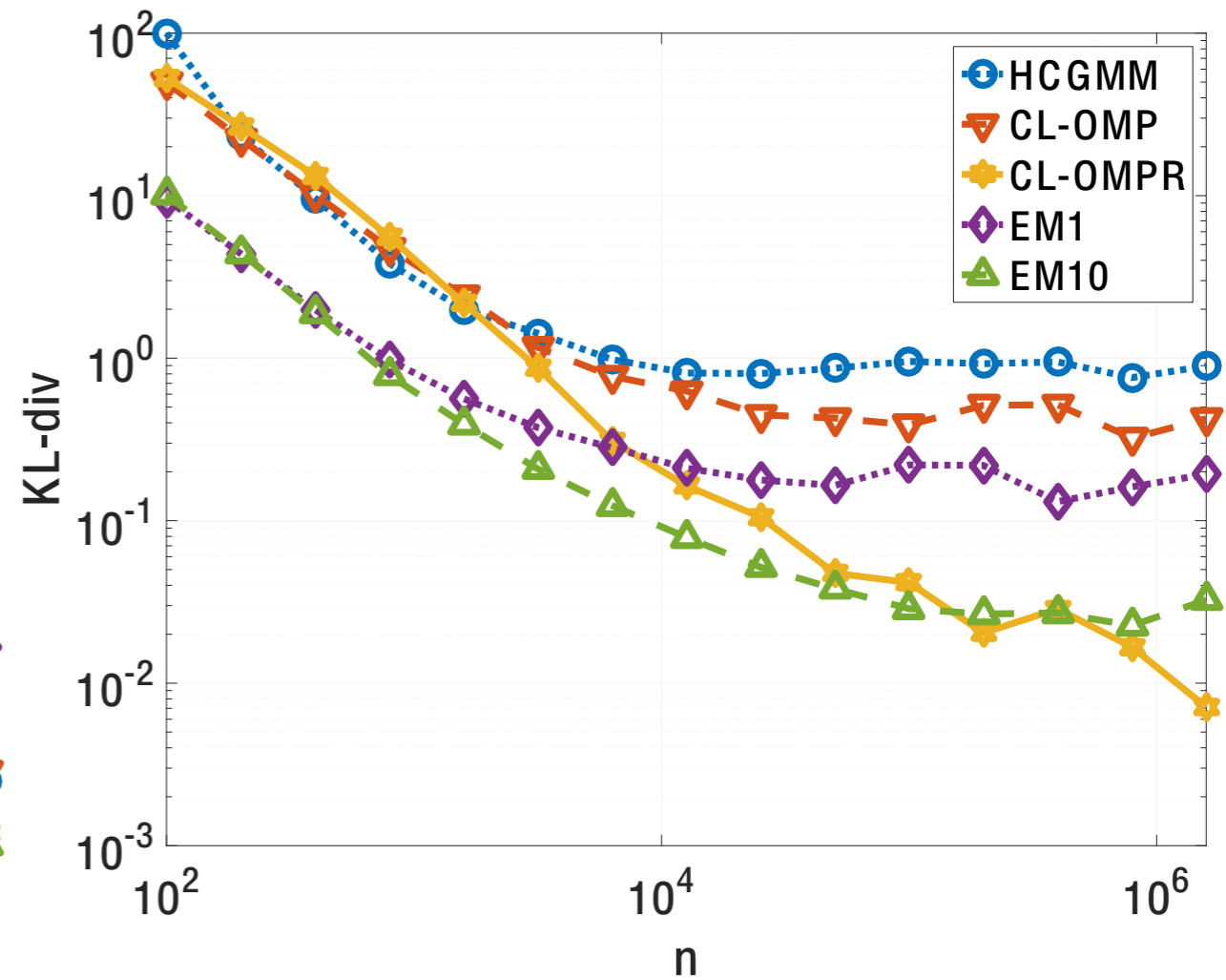
Sketching in practice

Results: $d = 10$

$K = 5$



$K = 20$



Sketching in practice

Learn from sketch:

- Depending on the task, find the suitable set:

$$\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$$

K-means: mixture of K diracs

GMM: mixture of K Gaussian

Generative model: distrib. parametrized by NN

- Solve the optimization problem:

Decoding problem:

- CL-OMP
- GD

$$\min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$

- Return the best parameter $\hat{\theta}$ and the distribution $\pi_{\hat{\theta}} \approx \pi$

| Sketching in practice

■ Another toy example:

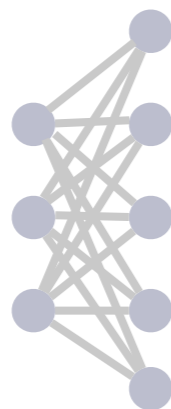
Generative model:

$$\mathcal{G}_\Theta = \left\{ \frac{1}{p} \sum_{j=1}^p \delta_{\text{NN}_\theta(\mathbf{z}_j)}; \theta \in \Theta \right\}$$



Latent space

$$\mathbf{z}_j \sim \Lambda$$



Sketching in practice

Another toy example:

Generative model:

$$\mathcal{G}_\Theta = \left\{ \frac{1}{p} \sum_{j=1}^p \delta_{\text{NN}_\theta(\mathbf{z}_j)}; \theta \in \Theta \right\}$$

Sketching op:

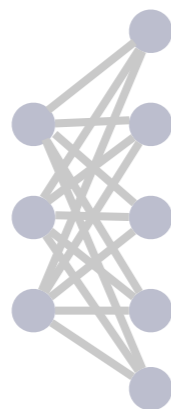
$$\mathcal{A}\pi_\theta = \frac{1}{p} \sum_{j=1}^p \Phi(\text{NN}_\theta(\mathbf{z}_j))$$

Optim.

SGD

Latent space

$$\mathbf{z}_j \sim \Lambda$$



Sketching in practice

Another toy example:

Generative model:

$$\mathcal{S}_\Theta = \left\{ \frac{1}{p} \sum_{j=1}^p \delta_{\text{NN}_\theta(\mathbf{z}_j)}; \theta \in \Theta \right\}$$

Sketching op:

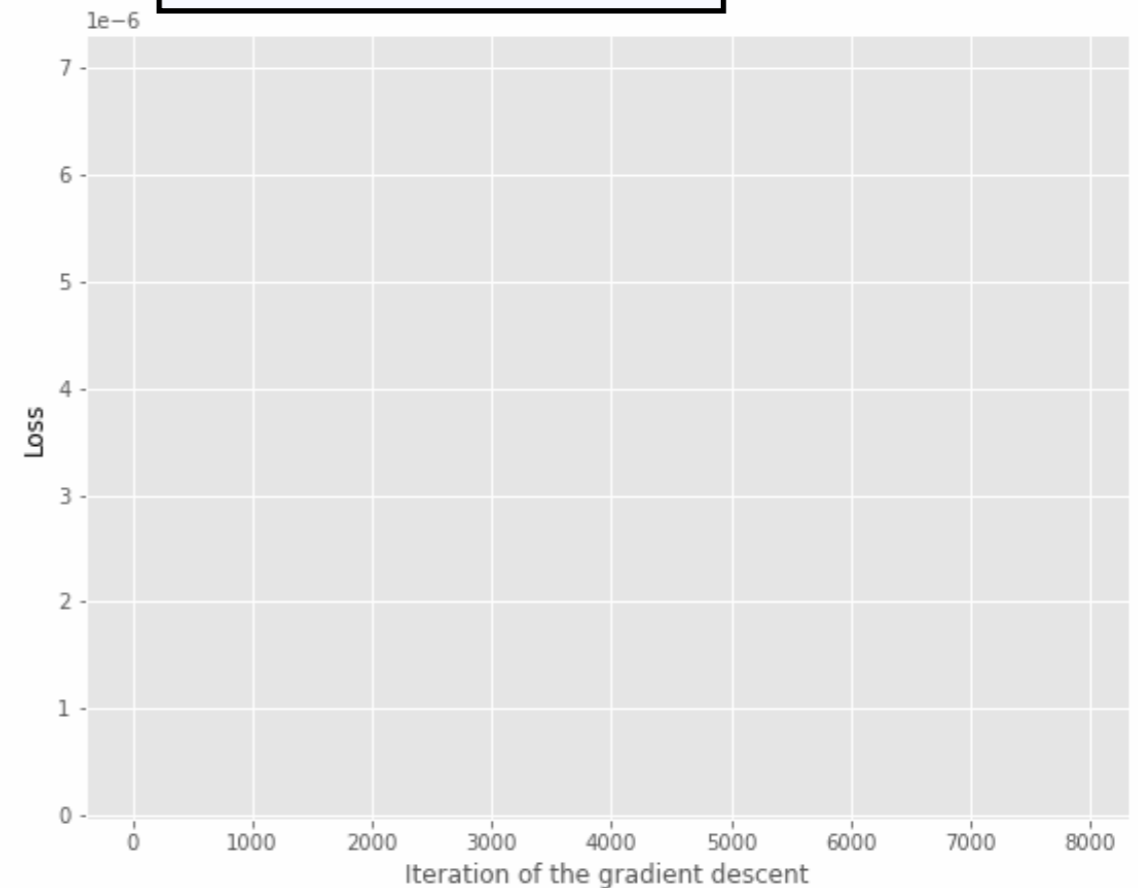
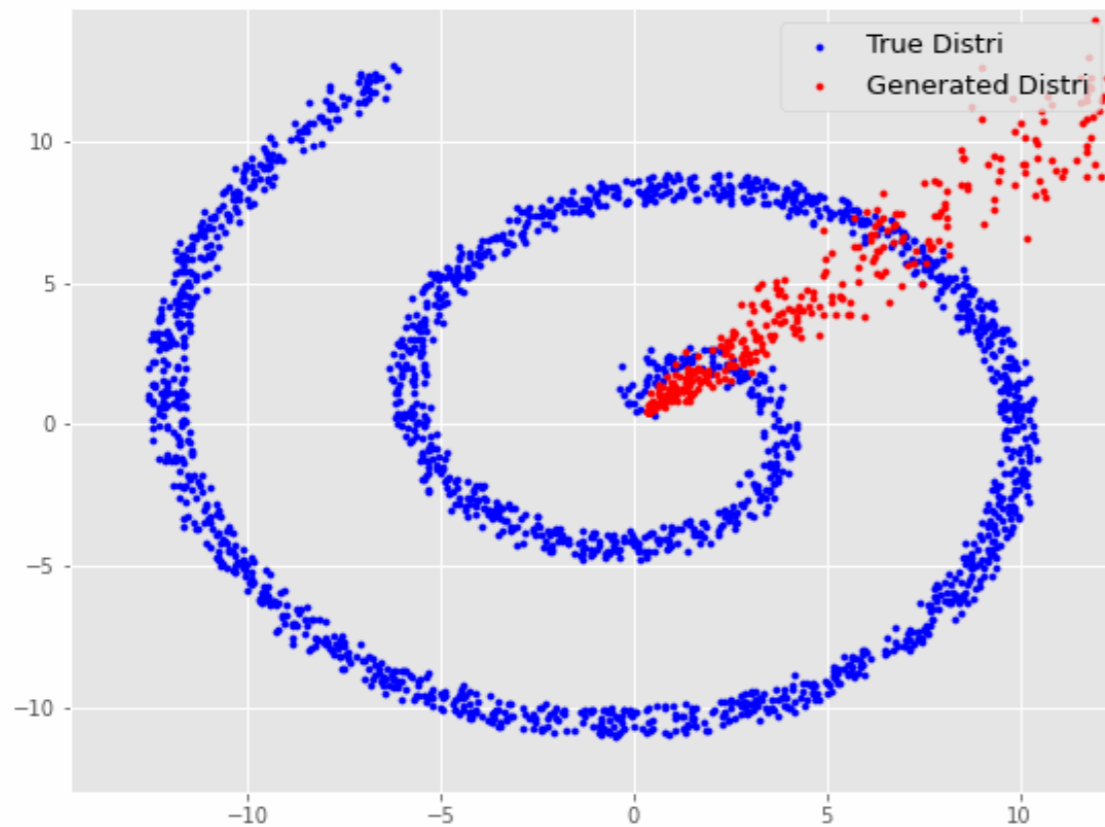
$$\mathcal{A}\pi_\theta = \frac{1}{p} \sum_{j=1}^p \Phi(\text{NN}_\theta(\mathbf{z}_j))$$

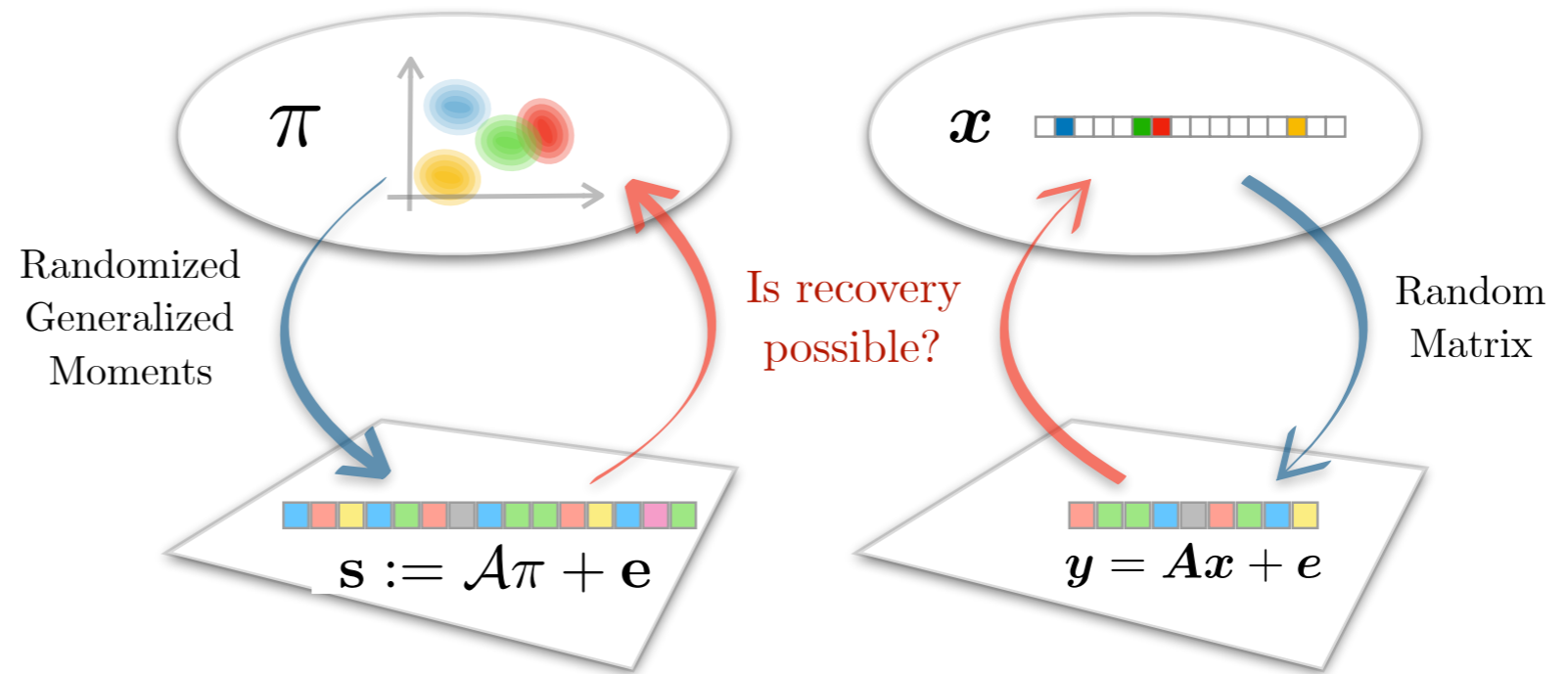
Optim.

SGD

Sketching for GAN

$$\min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$





Compressive Learning

- Theory of sketching
- Sketching in practice
- Theoretical guarantees
- Limitations & perspectives

| **Theoretical guarantees**

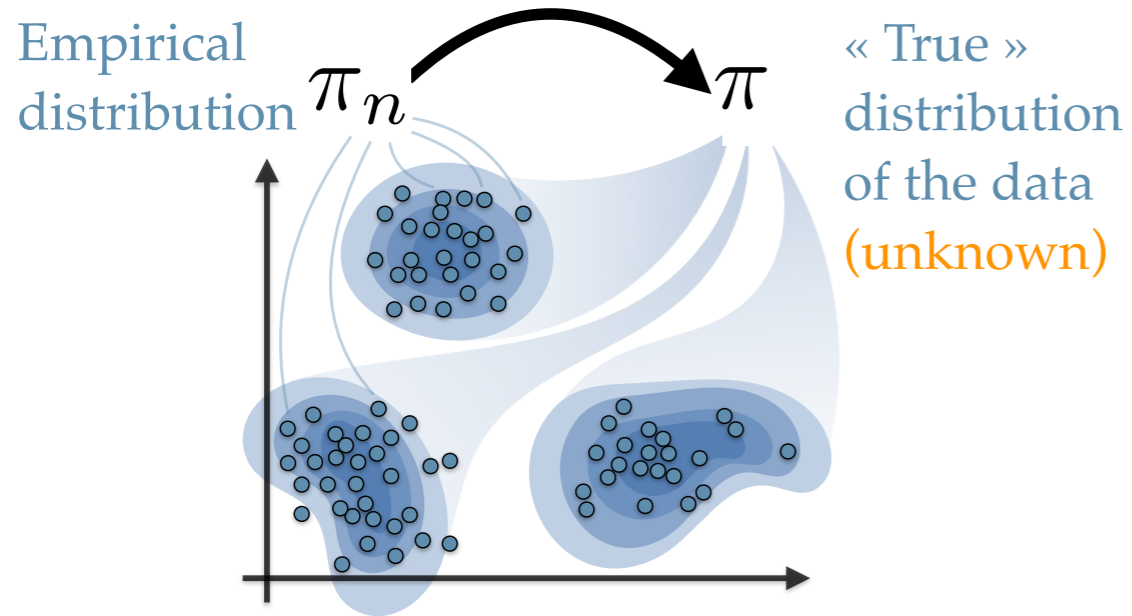
- **Analogy with compressed sensing**

Theoretical guarantees

Sketching operator

$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\Phi(\mathbf{x})]$$

Analogy with compressed sensing



We observe the sketch

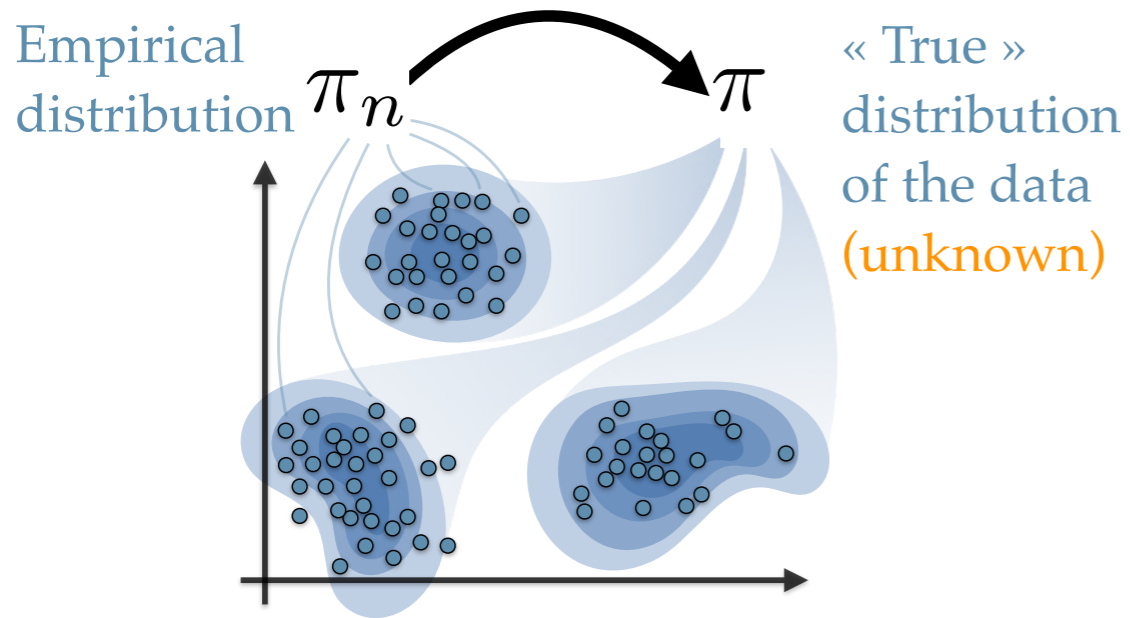
$$\mathbf{s} = \mathcal{A}\pi_n \in \mathbb{R}^m$$

Theoretical guarantees

Sketching operator

$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\Phi(\mathbf{x})]$$

Analogy with compressed sensing



We observe the sketch

$$\begin{aligned} \mathbf{s} &= \mathcal{A}\pi_n \in \mathbb{R}^m \\ &= \mathcal{A}\pi + \mathbf{e} \end{aligned}$$

where $\mathbf{e} := \mathcal{A}(\pi_n - \pi)$

noise

Noisy, linear and finite-dimensional measurements of a distribution

IN CS

Underdetermined $m < d$

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

IN Sketching

Very underdetermined $m < +\infty$

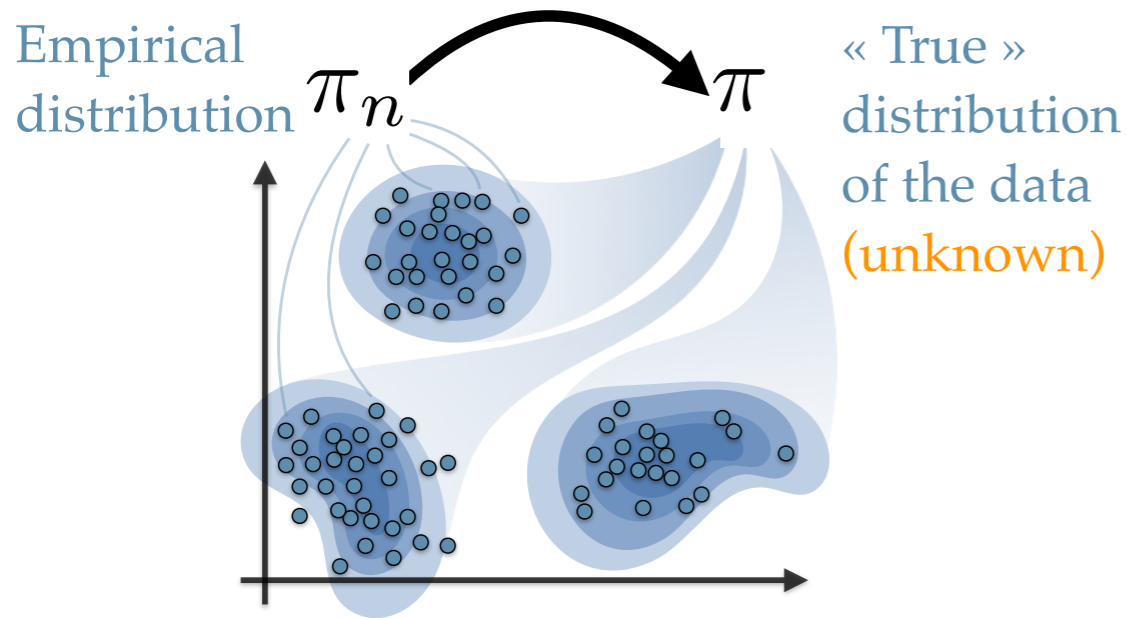
$$\mathbf{s} = \mathcal{A} \pi + \mathbf{e}$$

Theoretical guarantees

Sketching operator

$$\mathcal{A}\pi = \mathbb{E}_{\mathbf{x} \sim \pi} [\Phi(\mathbf{x})]$$

Analogy with compressed sensing



We observe the sketch

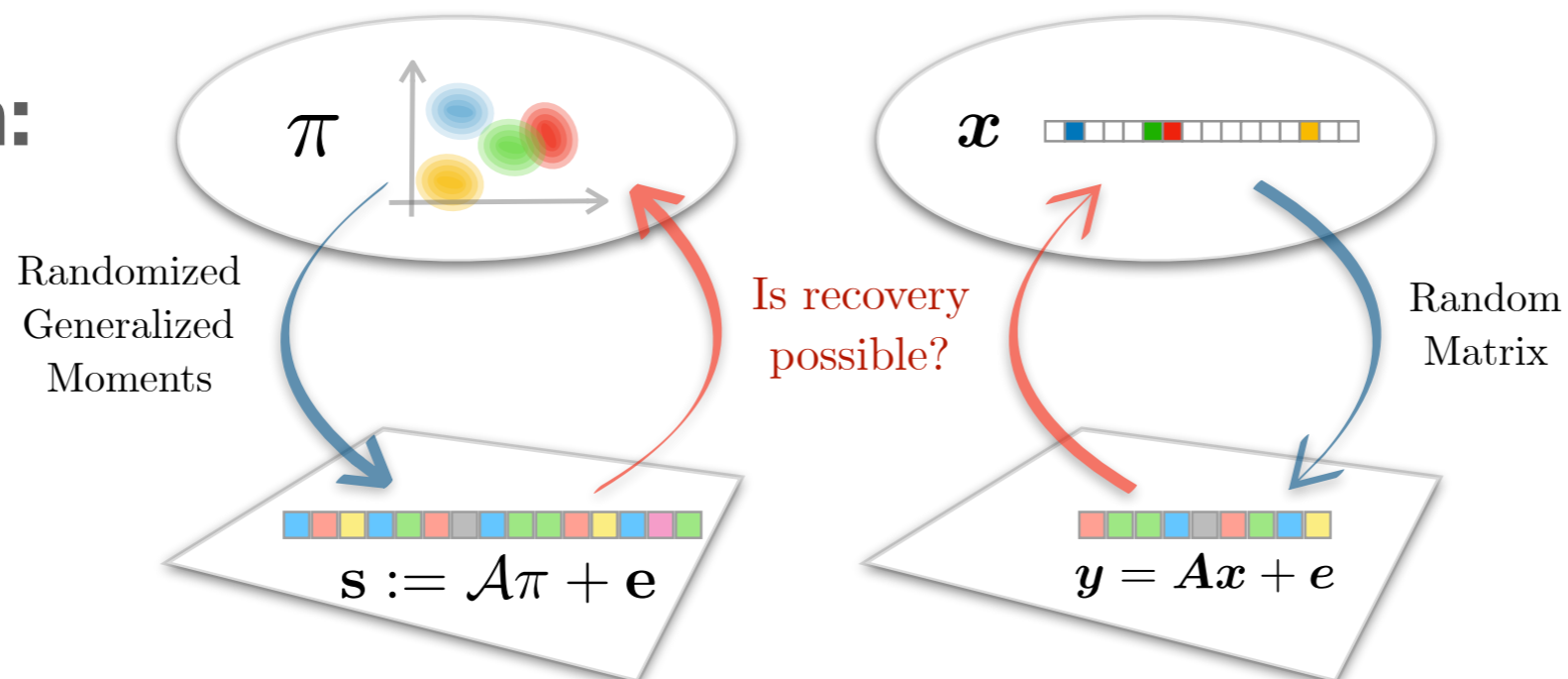
$$\begin{aligned} \mathbf{s} &= \mathcal{A}\pi_n \in \mathbb{R}^m \\ &= \mathcal{A}\pi + \mathbf{e} \end{aligned}$$

where $\mathbf{e} := \mathcal{A}(\pi_n - \pi)$

noise

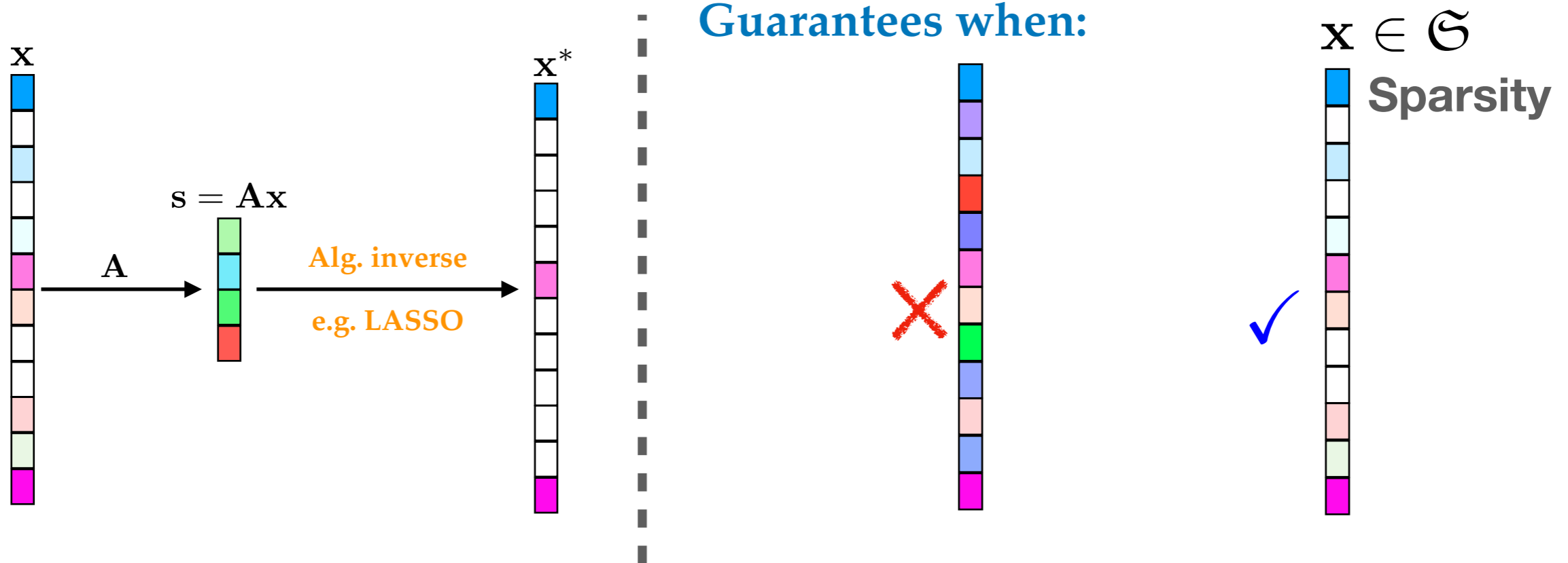
Noisy, linear and finite-dimensional measurements of a distribution

Question:

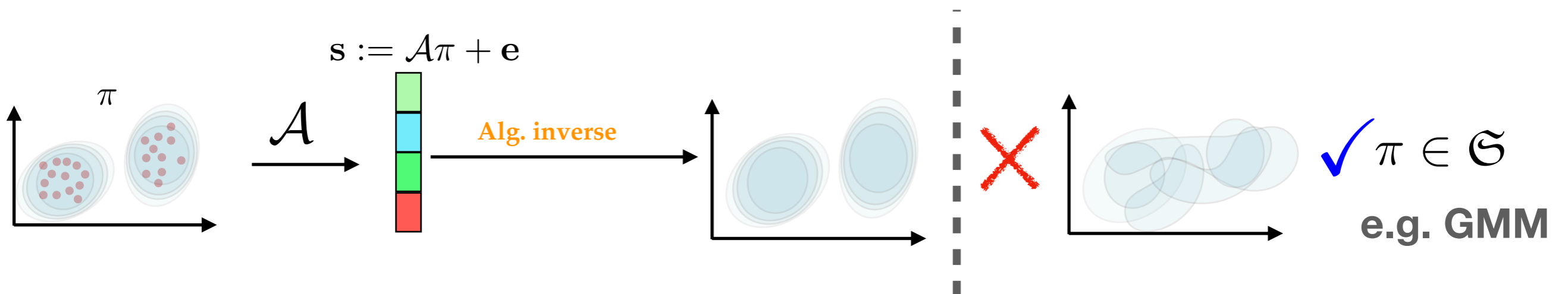


Theoretical guarantees

Analogy with compressed sensing



...Need a « low-dimensional » distribution



Theoretical guarantees

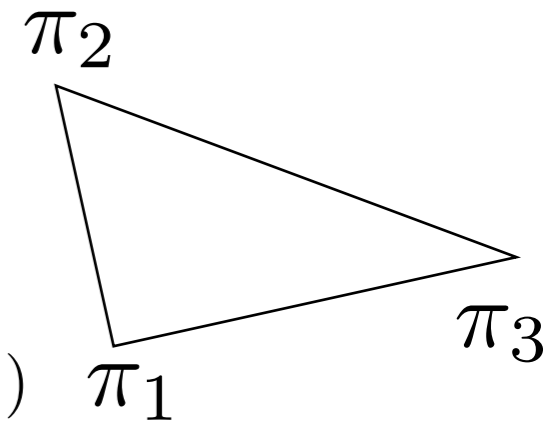
■ Metric between distributions:

■ Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**

1) $D(\pi, \pi') = 0 \iff \pi = \pi'$ 3) $D(\pi_1, \pi_2) \leq D(\pi_1, \pi_3) + D(\pi_3, \pi_2)$

2) $D(\pi, \pi') = D(\pi', \pi)$

Quantifies the distance between distrib.



Theoretical guarantees

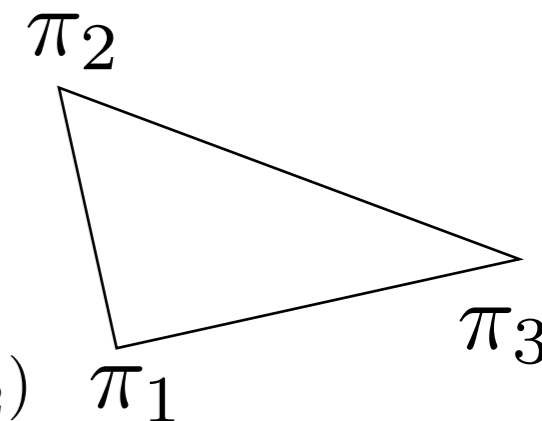
■ Metric between distributions:

■ Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**

1) $D(\pi, \pi') = 0 \iff \pi = \pi'$ 3) $D(\pi_1, \pi_2) \leq D(\pi_1, \pi_3) + D(\pi_3, \pi_2)$

2) $D(\pi, \pi') = D(\pi', \pi)$

Quantifies the distance between distrib.



■ **Numerous examples** in the literature...

Total Variation

$$D(\pi, \pi') = \|\pi - \pi'\|_{\text{TV}}$$

if densities f, g : $= \frac{1}{2} \|f - g\|_{L_1}$

Theoretical guarantees

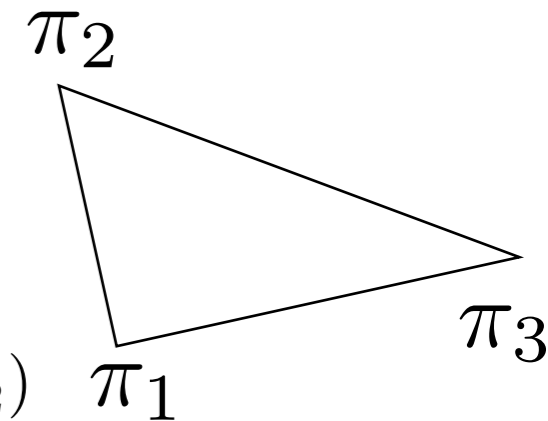
Metric between distributions:

Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**

1) $D(\pi, \pi') = 0 \iff \pi = \pi'$ 3) $D(\pi_1, \pi_2) \leq D(\pi_1, \pi_3) + D(\pi_3, \pi_2)$

2) $D(\pi, \pi') = D(\pi', \pi)$

Quantifies the distance between distrib.



Numerous examples in the literature...

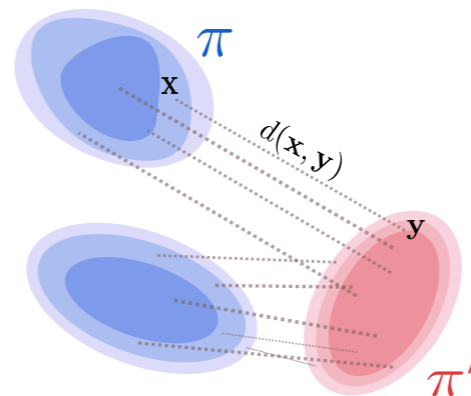
Total Variation

$$D(\pi, \pi') = \|\pi - \pi'\|_{\text{TV}}$$

if densities f, g : $= \frac{1}{2} \|f - g\|_{L_1}$

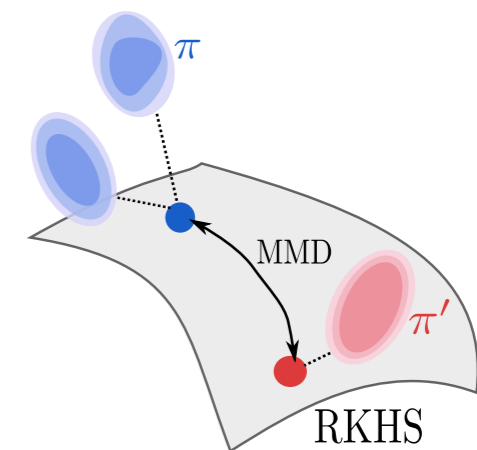
Optimal Transport

$$D(\pi, \pi') = W_p(\pi, \pi')$$



MMD

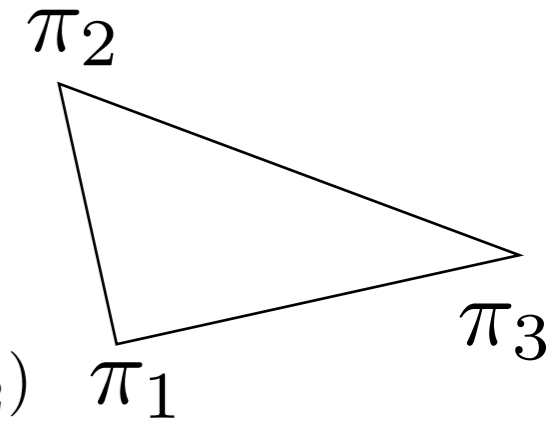
$$D(\pi, \pi') = \|\pi - \pi'\|_{\kappa}$$



Theoretical guarantees

Metric between distributions:

Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**



1) $D(\pi, \pi') = 0 \iff \pi = \pi'$ 3) $D(\pi_1, \pi_2) \leq D(\pi_1, \pi_3) + D(\pi_3, \pi_2)$

2) $D(\pi, \pi') = D(\pi', \pi)$

Quantifies the distance between distrib.

Numerous examples in the literature...

Total Variation

$$D(\pi, \pi') = \|\pi - \pi'\|_{\text{TV}}$$

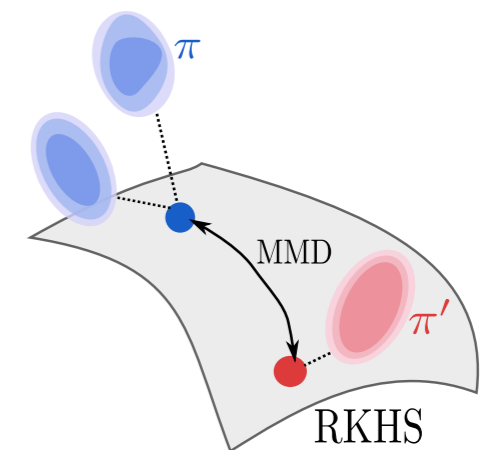
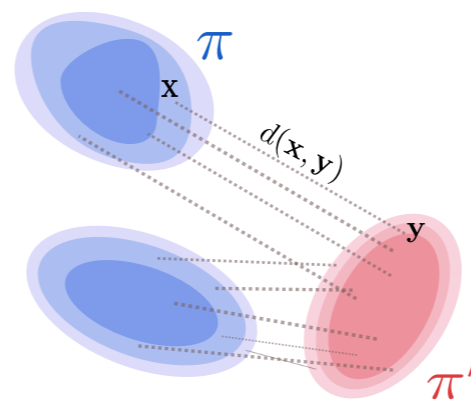
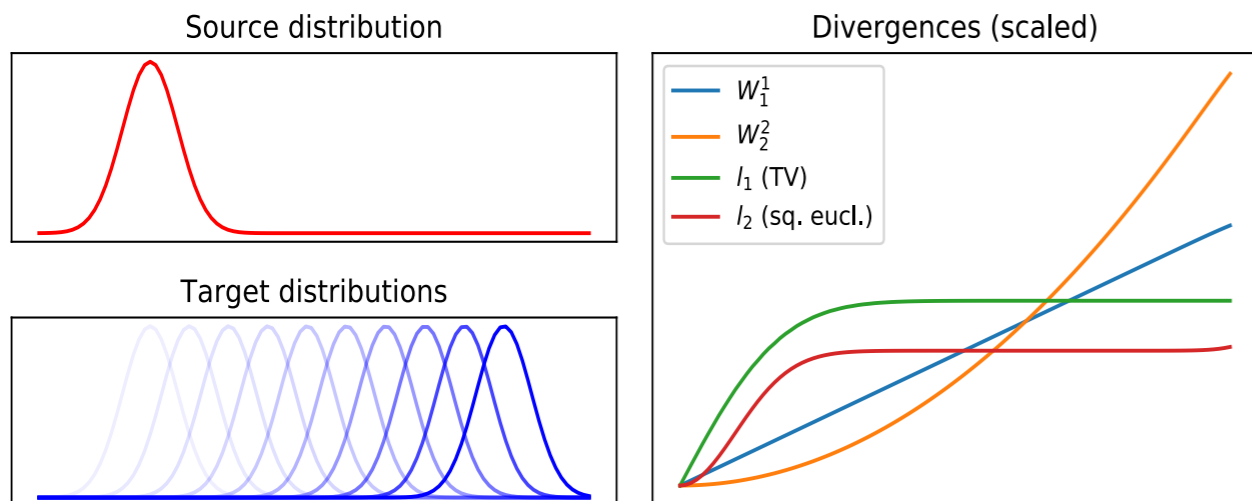
if densities f, g : $= \frac{1}{2} \|f - g\|_{L_1}$

Optimal Transport

$$D(\pi, \pi') = W_p(\pi, \pi')$$

MMD

$$D(\pi, \pi') = \|\pi - \pi'\|_{\kappa}$$



Theoretical guarantees

■ The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Theoretical guarantees

The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

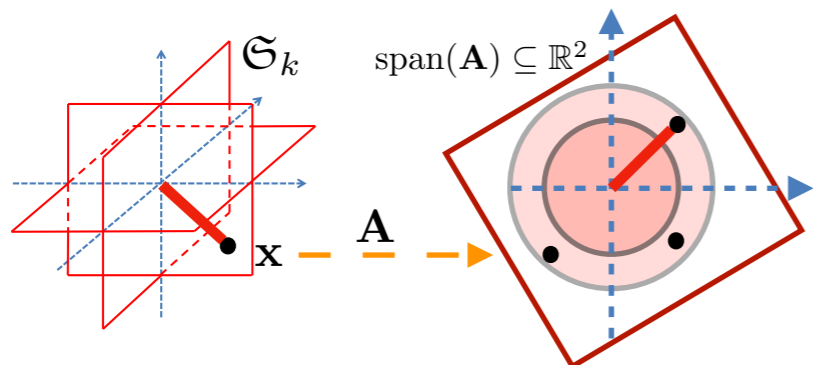
$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Connection with the RIP:

$$\begin{aligned} \exists \delta_k \in [0, 1[\quad \forall \mathbf{x} \text{ } k\text{-sparse} \\ (1 - \delta_k) \|\mathbf{x}\|_2^2 \leq \|\mathbf{A}\mathbf{x}\|_2^2 \leq (1 + \delta_k) \|\mathbf{x}\|_2^2 \end{aligned}$$

E.g. $D(\pi, \pi') = \|\pi - \pi'\|_{\text{TV}}$

$$\begin{aligned} \exists C > 0 \quad \forall x = \pi_\theta - \pi_{\theta'} \\ \|x\|_{\text{TV}} \leq C \|\mathcal{A}x\|_2 \end{aligned}$$



same kind of idea

Theoretical guarantees

■ The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Suppose **Lower LRIP** holds:

Take $\pi \in \mathcal{P}(\mathbb{R}^d)$ emp. distrib. π_n

Sketch $\mathbf{s} = \mathcal{A}\pi_n \in \mathbb{R}^m$

Solve:

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$

Theoretical guarantees

■ The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathfrak{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Suppose **Lower LRIP** holds:

Take $\pi \in \mathcal{P}(\mathbb{R}^d)$ emp. distrib. π_n

Sketch $\mathbf{s} = \mathcal{A}\pi_n \in \mathbb{R}^m$

Solve:

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$

$$\implies D(\pi_{\hat{\theta}}, \pi) \leq d^\circ(\pi, \mathfrak{S}_\Theta) + 2C \|\mathcal{A}\pi - \mathcal{A}\pi_n\|_2$$

Theoretical guarantees

The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Suppose **Lower LRIP** holds:

Take $\pi \in \mathcal{P}(\mathbb{R}^d)$ emp. distrib. π_n

Sketch $\mathbf{s} = \mathcal{A}\pi_n \in \mathbb{R}^m$

Solve:

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$

$$\implies D(\pi_{\hat{\theta}}, \pi) \leq d^\circ(\pi, \mathcal{S}_\Theta) + 2C \|\mathcal{A}\pi - \mathcal{A}\pi_n\|_2$$

distance between the true distrib. and the estimated one

Theoretical guarantees

■ The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Suppose **Lower LRIP** holds:

Take $\pi \in \mathcal{P}(\mathbb{R}^d)$ emp. distrib. π_n

Sketch $\mathbf{s} = \mathcal{A}\pi_n \in \mathbb{R}^m$

Solve:

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$

$$\implies D(\pi_{\hat{\theta}}, \pi) \leq d^\circ(\pi, \mathcal{S}_\Theta)$$

How far is the true distrib. from the model = approx. error

Theoretical guarantees

■ The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathfrak{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Suppose **Lower LRIP** holds:

Take $\pi \in \mathcal{P}(\mathbb{R}^d)$ emp. distrib. π_n

Sketch $\mathbf{s} = \mathcal{A}\pi_n \in \mathbb{R}^m$

Solve:

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$

$$\implies D(\pi_{\hat{\theta}}, \pi) \leq d^\circ(\pi, \mathfrak{S}_\Theta) + 2C \|\mathcal{A}\pi - \mathcal{A}\pi_n\|_2$$

Error term $\mathcal{O}(n^{-1/2})$

Theoretical guarantees

■ The lower RIP:

■ Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**

■ The **model set** related to the task $\mathfrak{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

Suppose **Lower LRIP** holds:

Take $\pi \in \mathcal{P}(\mathbb{R}^d)$ emp. distrib. π_n

Sketch $\mathbf{s} = \mathcal{A}\pi_n \in \mathbb{R}^m$

Solve:

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \|\mathbf{s} - \mathcal{A}\pi_\theta\|_2$$

$$\implies D(\pi_{\hat{\theta}}, \pi) \leq d^\circ(\pi, \mathfrak{S}_\Theta) + 2C \|\mathcal{A}\pi - \mathcal{A}\pi_n\|_2$$

$$\implies \text{If } \pi \in \mathfrak{S}_\Theta \text{ then } D(\pi_{\hat{\theta}}, \pi) \xrightarrow{n \rightarrow +\infty} 0$$

Theoretical guarantees

■ The lower RIP:

- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

■ We want the LRIP:

- The Lower RIP implies interesting theoretical guarantees
- However difficult to prove in general

Theoretical guarantees

The lower RIP:

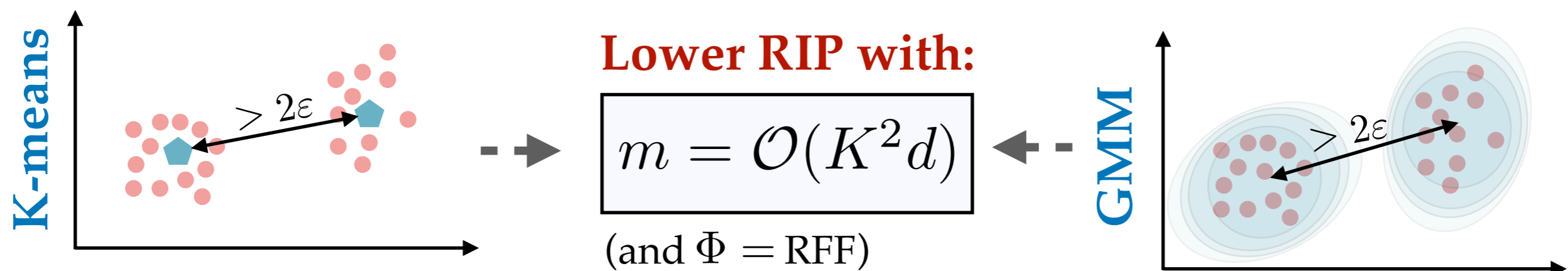
- Let $D : \mathcal{P}(\mathbb{R}^d) \times \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}_+$ be a **metric**
- The **model set** related to the task $\mathcal{S}_\Theta = \{\pi_\theta; \theta \in \Theta\}$

Lower RIP

$$\forall \theta, \theta' \in \Theta, D(\pi_\theta, \pi_{\theta'}) \leq C \|\mathcal{A}\pi_\theta - \mathcal{A}\pi_{\theta'}\|_2$$

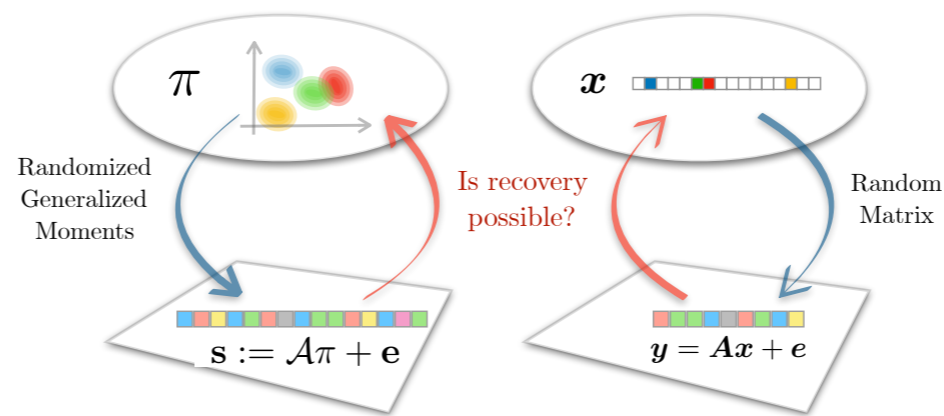
We want the LRIP:

- The Lower RIP implies interesting theoretical guarantees
- However difficult to prove in general **with separability of the clusters**

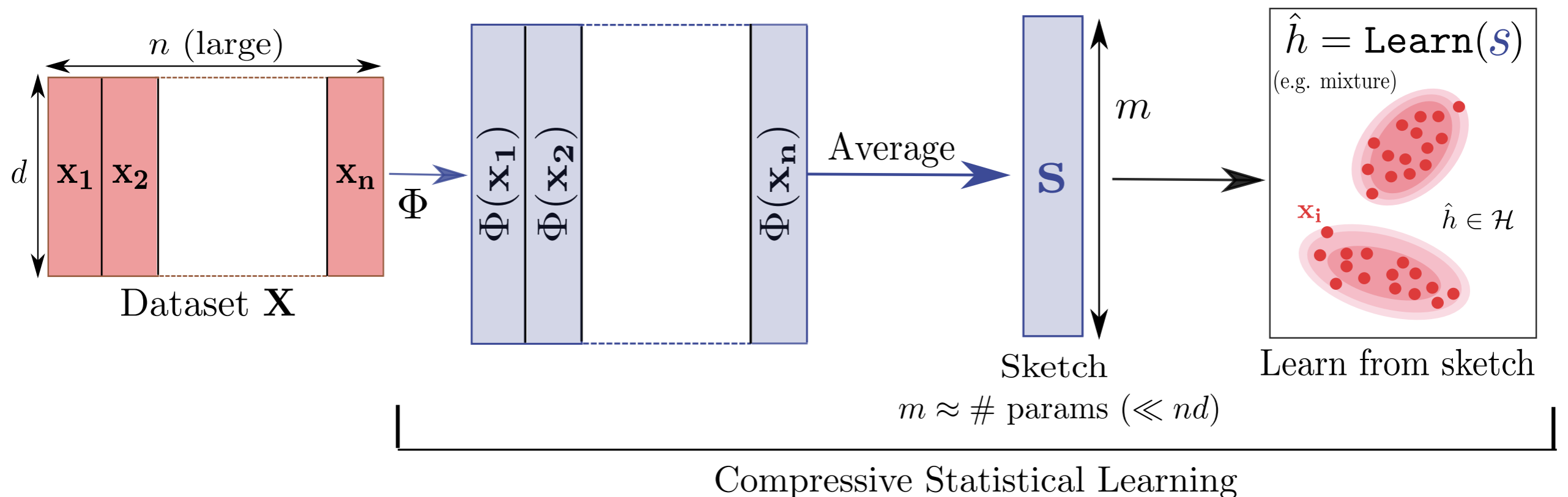


Conclusion

Sketching theory

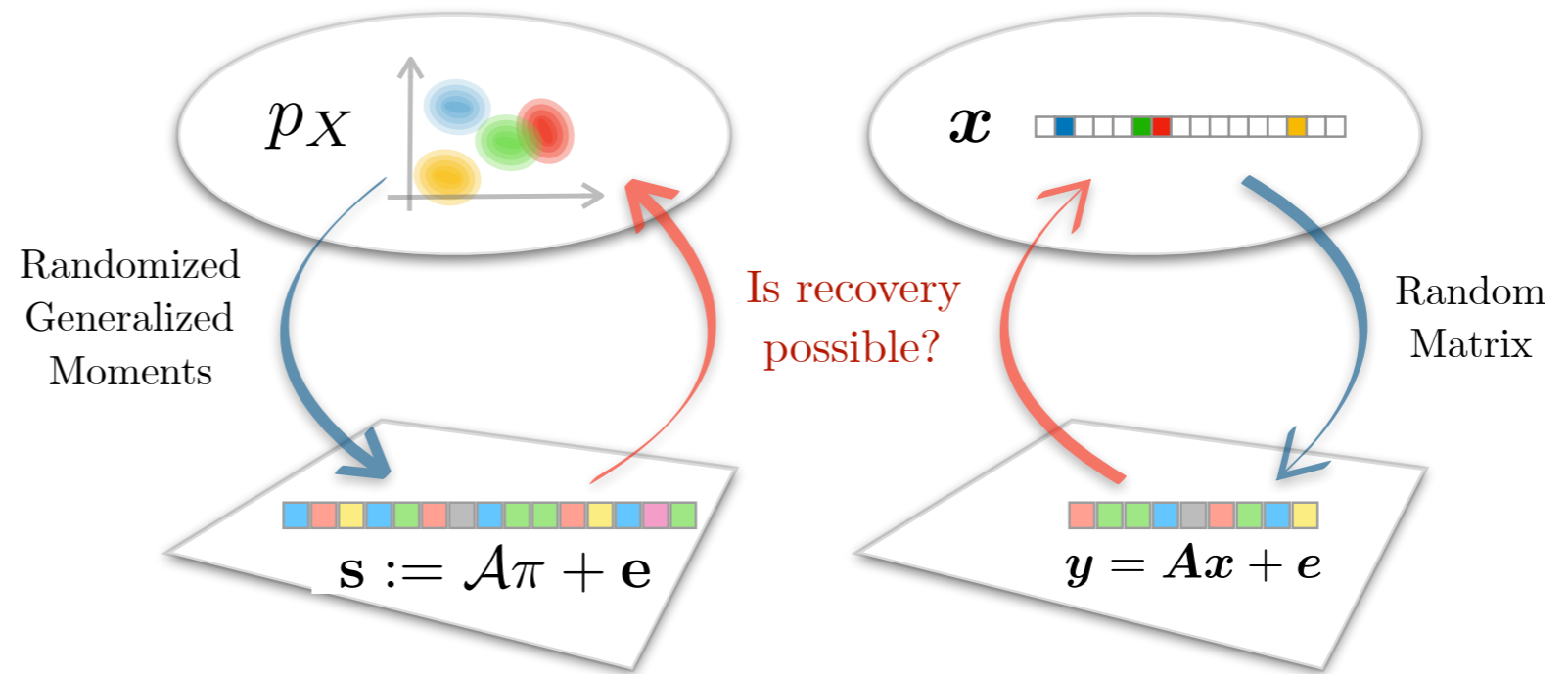


Method for **resource efficient large-scale machine learning**



Comes with statistical guarantees inspired by **compressed sensing**

Useful for **online, distributed and private** learning



Compressive Learning

- Theory of sketching
- Sketching in practice
- Theoretical guarantees
- Limitations & perspectives

| Limitations & perspectives

■ Complexity in space

■ Sketching reduces drastically the dimension but....

Need to store somewhere:

$$\mathbf{W} \in \mathbb{R}^{m \times d}$$

dense random matrix ...

Need to calculate the sketch: $\mathcal{O}(nmd)$

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \exp(-i \mathbf{W} \mathbf{x}_i)$$

Limitations & perspectives

Complexity in space

- Sketching reduces drastically the dimension but...

Need to store somewhere:

$$\mathbf{W} \in \mathbb{R}^{m \times d}$$

dense random matrix ...

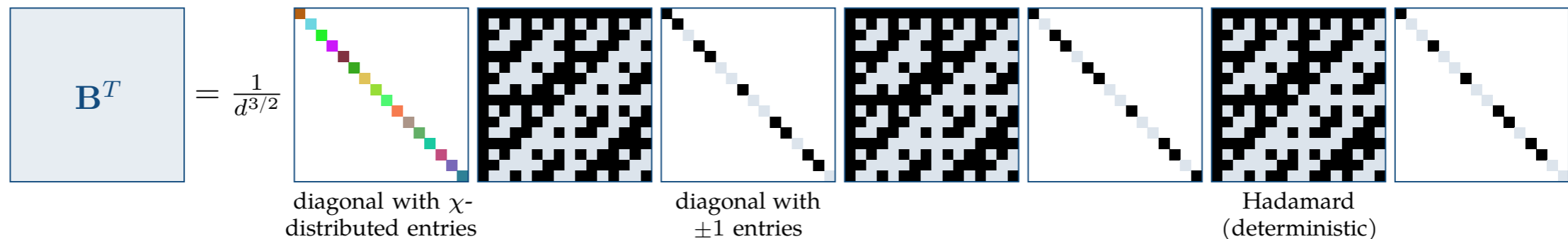
Need to calculate the sketch: $\mathcal{O}(nmd)$

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \exp(-i \mathbf{W} \mathbf{x}_i)$$

Perspectives

- Use random structured matrices:

$$\mathbf{W} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \dots & \mathbf{B}_b \end{bmatrix} \quad d_p = 2^q = d$$



Limitations & perspectives

Complexity in space

Sketching reduces drastically the dimension but....

Need to store somewhere:

$$\mathbf{W} \in \mathbb{R}^{m \times d}$$

dense random matrix ...

Need to calculate the sketch: $\mathcal{O}(nmd)$

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \exp(-i\mathbf{W}\mathbf{x}_i)$$

Fast transform

$$\mathcal{O}(nm \ln(d))$$

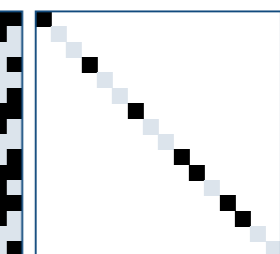
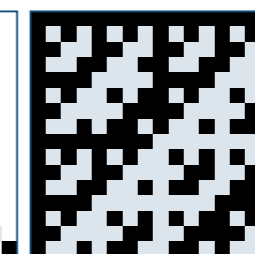
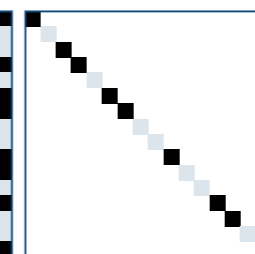
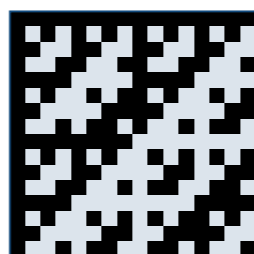
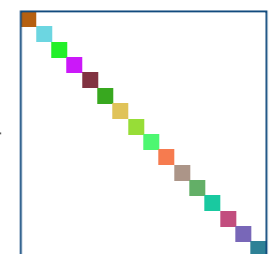
Perspectives

Use random structured matrices:

sparse

$$\mathbf{W} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \dots & \mathbf{B}_b \end{bmatrix} \quad d_p = 2^q = d$$

$$\mathbf{B}^T = \frac{1}{d^{3/2}}$$



Limitations & perspectives

Complexity in space

Sketching reduces drastically the dimension but....

Need to store somewhere:

$$\mathbf{W} \in \mathbb{R}^{m \times d}$$

dense random matrix ...

Need to calculate the sketch: $\mathcal{O}(nmd)$

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \exp(-i\mathbf{W}\mathbf{x}_i)$$

Fast transform

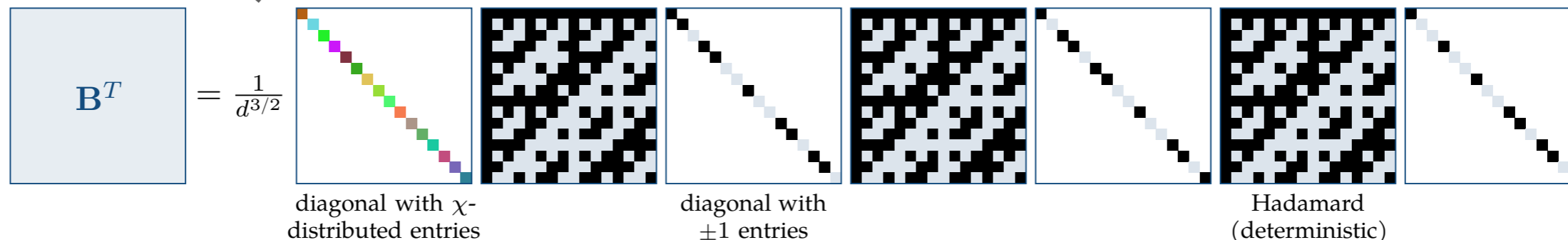
$$\mathcal{O}(nm \ln(d))$$

Perspectives

Use random structured matrices + parallelize

sparse

$$\mathbf{W} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \dots & \mathbf{B}_b \end{bmatrix} \quad d_p = 2^q = d$$



Limitations & perspectives

Complexity in space

Sketching reduces drastically the dimension but....

Need to store somewhere:

$$\mathbf{W} \in \mathbb{R}^{m \times d}$$

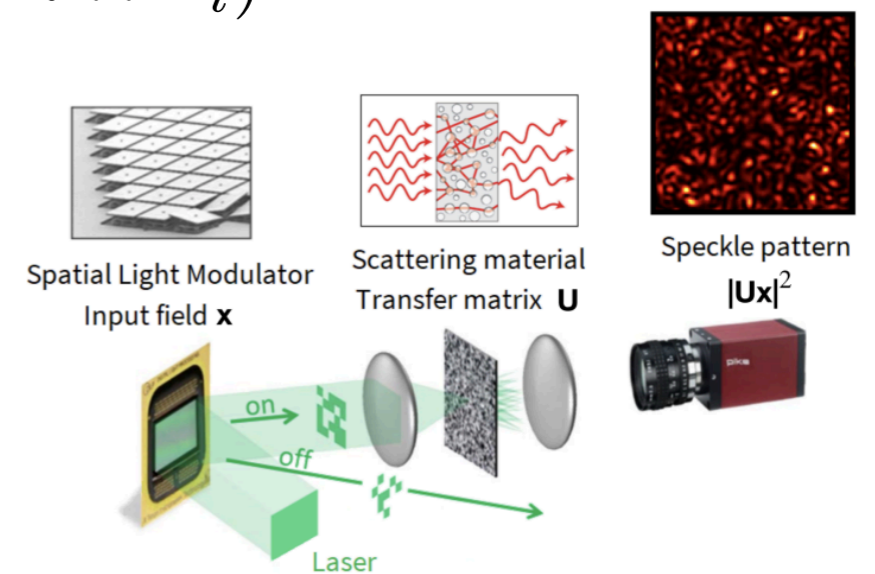
Need to **calculate the sketch**: $\mathcal{O}(nmd)$

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \exp(-i\mathbf{W}\mathbf{x}_i)$$

Perspectives

Use random structured matrices

Optical Processing Unit (OPU)

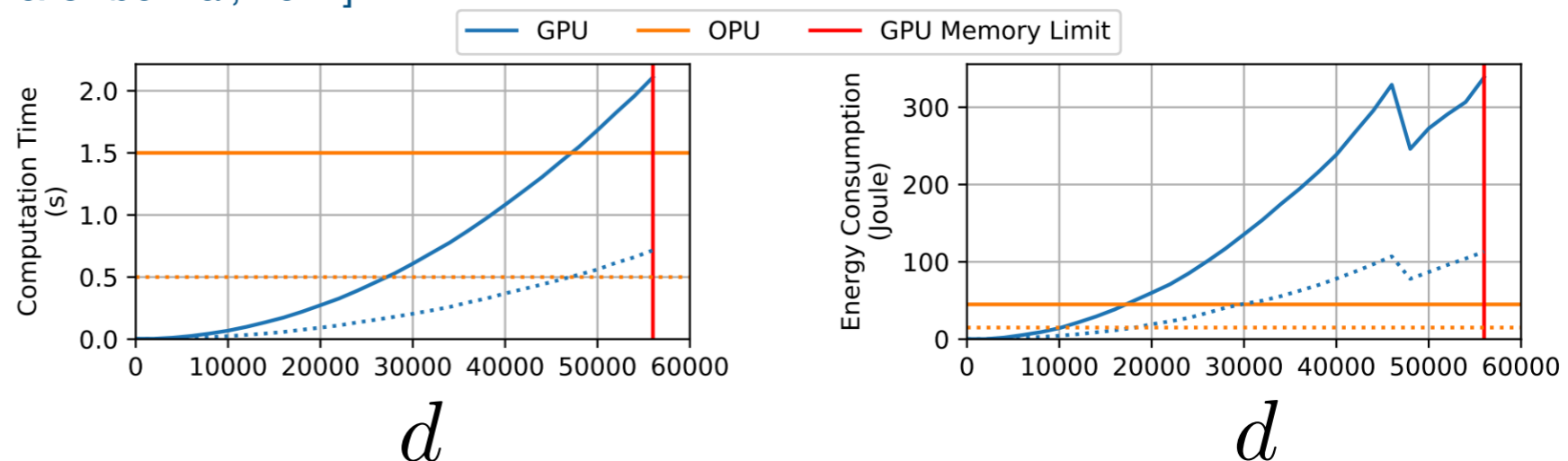


[Giffon & Gribonval, 2022]

- Constant time for matrix multiplication (with random matrix)

$$\mathbf{W}\mathbf{x}_i \text{ in } \mathcal{O}(1)$$

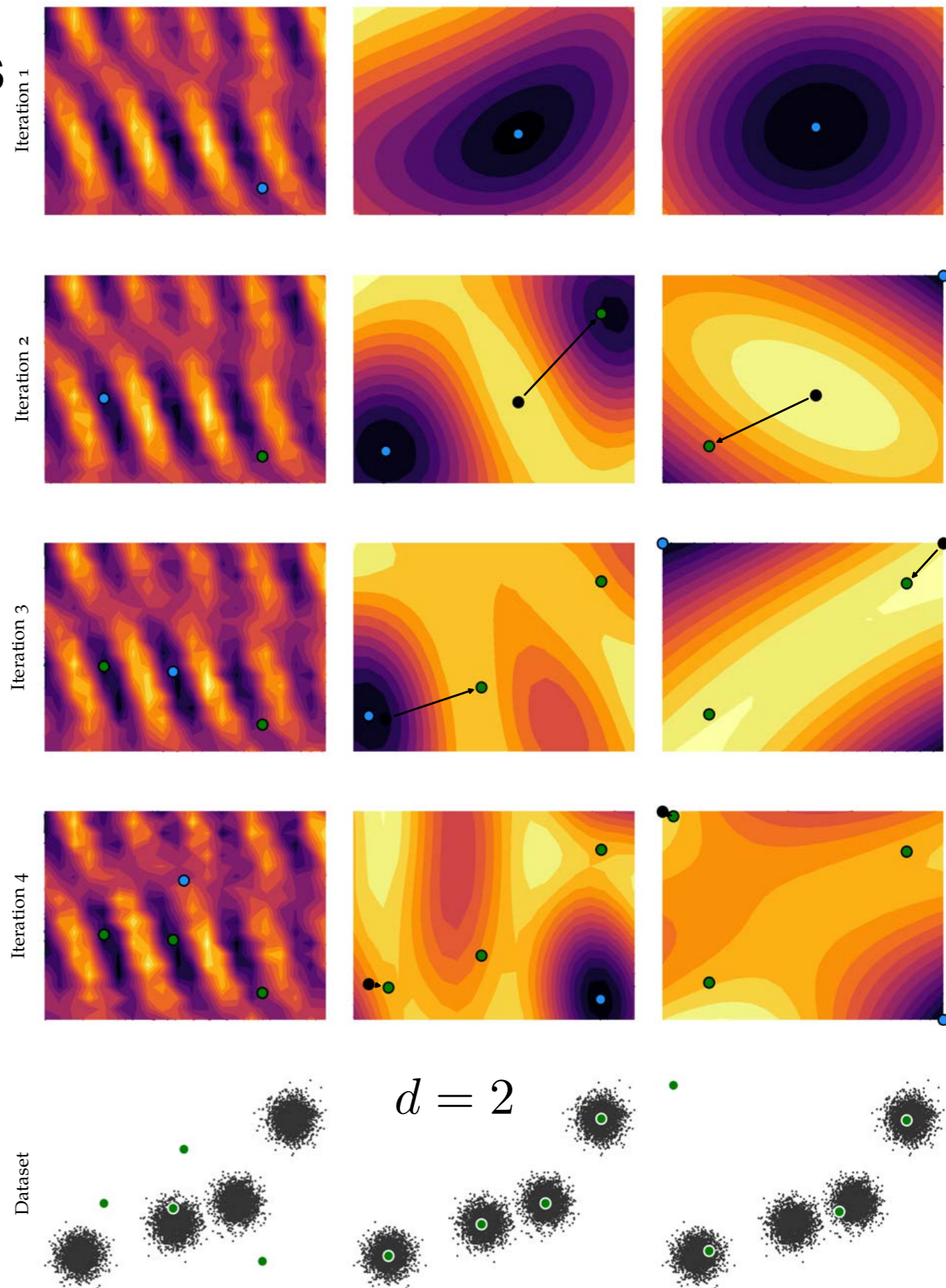
- Constant power consumption



Limitations & perspectives

Hyperparameter

$$W_{ij} \sim \mathcal{N}(0, \sigma^2) \quad \text{Very sensitive}$$



Limitations & perspectives

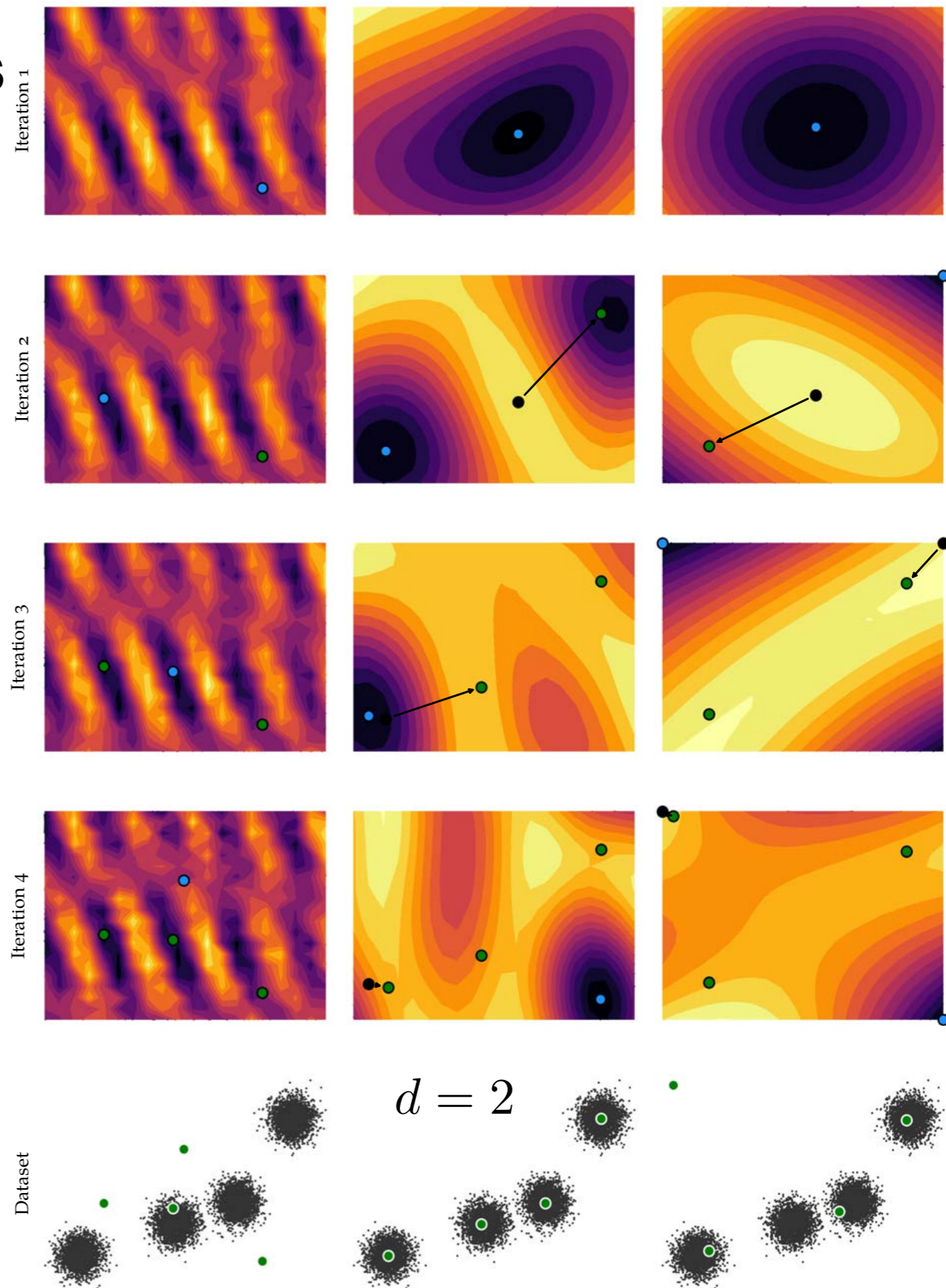
Hyperparameter

$$W_{ij} \sim \mathcal{N}(0, \sigma^2) \quad \text{Very sensitive}$$

Theoretical guarantees

Lower RIP difficult to prove

Decoding = non-convex problem



Limitations & perspectives

Hyperparameter

$$W_{ij} \sim \mathcal{N}(0, \sigma^2) \quad \text{Very sensitive}$$

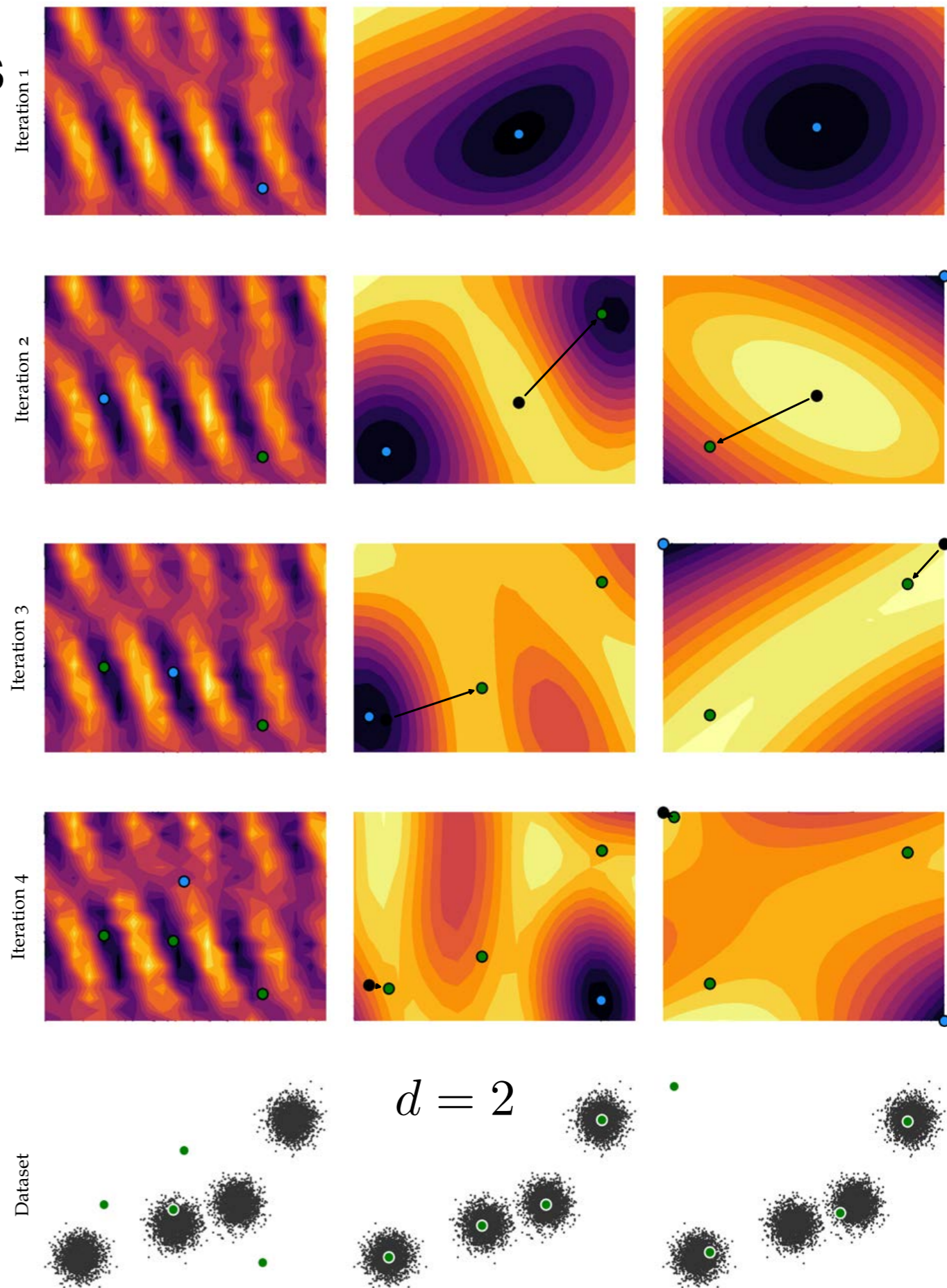
Theoretical guarantees

Lower RIP difficult to prove

Decoding = non-convex problem

More tasks ?

Supervised machine learning ?



Limitations & perspectives

Hyperparameter

$$W_{ij} \sim \mathcal{N}(0, \sigma^2) \quad \text{Very sensitive}$$

Theoretical guarantees

Lower RIP difficult to prove

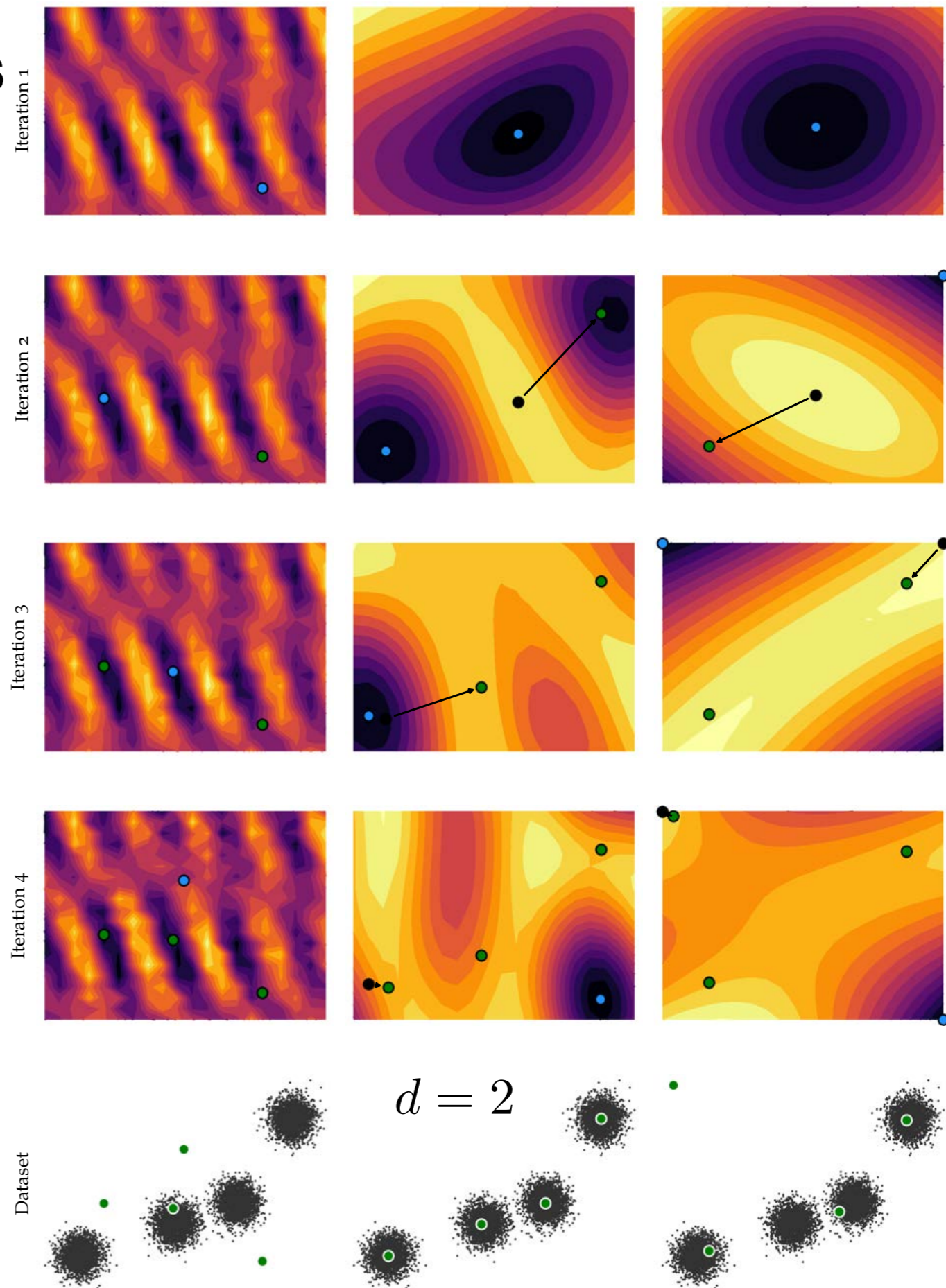
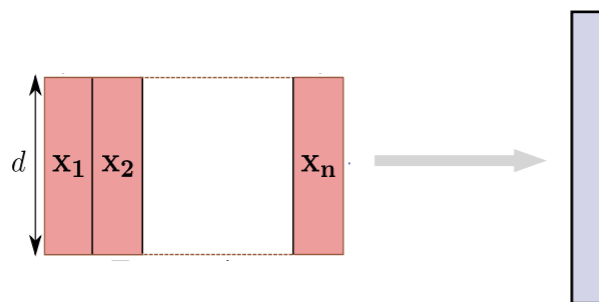
Decoding = non-convex problem

More tasks ?

Supervised machine learning ?

Privacy

[Chatalic, 2020]



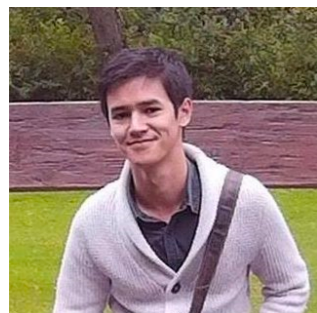
The sketching theory



- Series of past work with **Rémi Gribonval**
- And also Anthony Bourrier, **Nicolas Keriven**, **Antoine Chatalic**, Ayoub Belhadji, Luc Giffon, Gilles Puy, Nicolas Tremblay, Yann Traonmilin, Clément Elvira, Patrick Perez, Mike Davies, Gilles Blanchard, Pierre Vandergheynst, Laurent Jacques, **Vincent Schellekens**, Florimond Houssiau, Phil Schniter, Evan Byrne, ...

Sketching for large-scale learning of mixture models

Nicolas Keriven



Efficient and privacy-preserving compressive learning

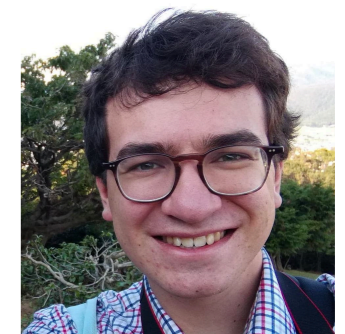
Antoine Chatalic



also thanks to
Rémi Flamary
for some figures...

Extending the Compressive Statistical
Learning Framework:
Quantization, Privacy, and Beyond

Vincent Schellekens



IEEE Signal Processing MAGAZINE

Volume 38 | Number 5 | September 2021

SKETCHING DATA SETS FOR LARGE-SCALE LEARNING

Keeping Only What You Need

Diagnosis/Prognosis
of COVID-19 Chest Images

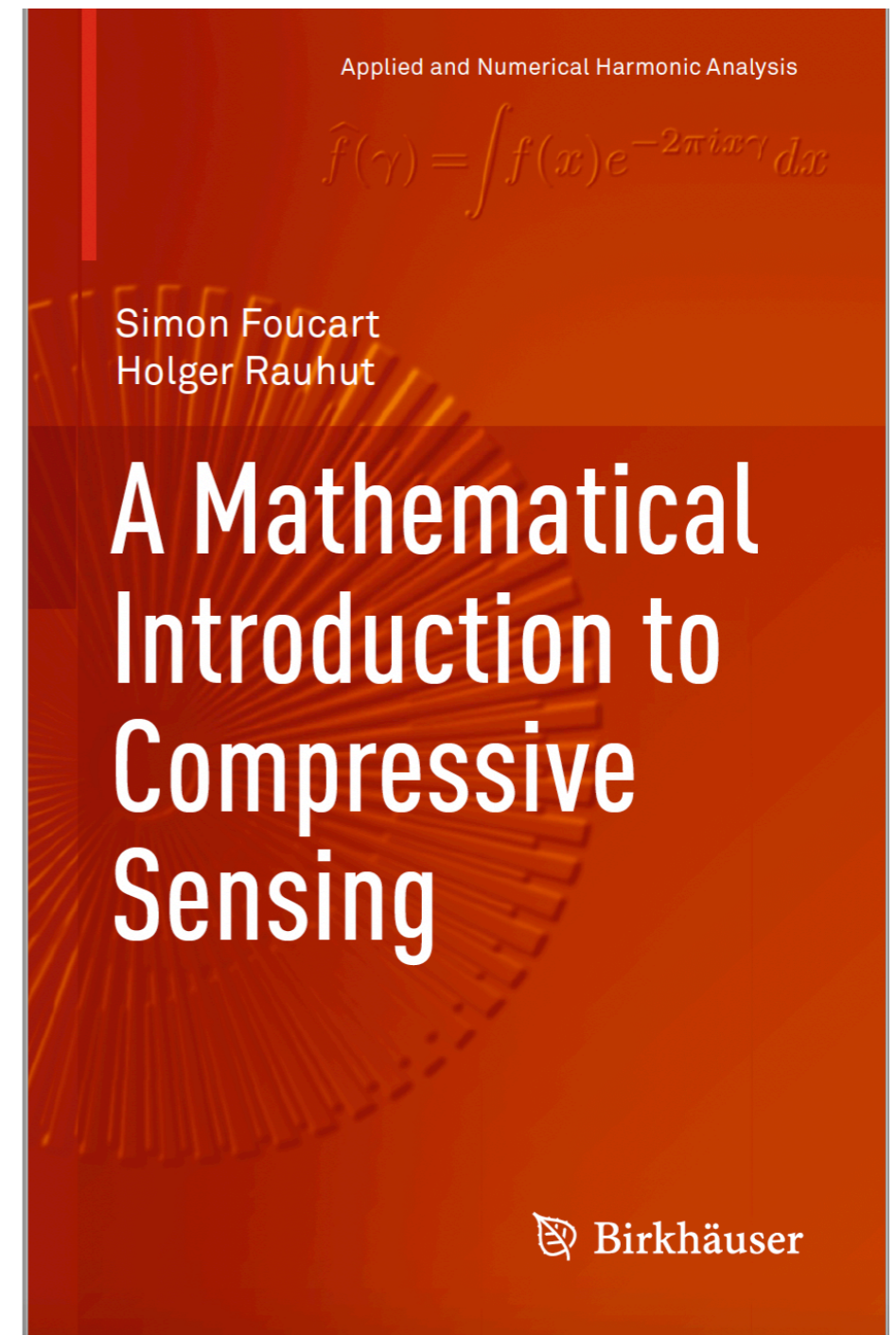
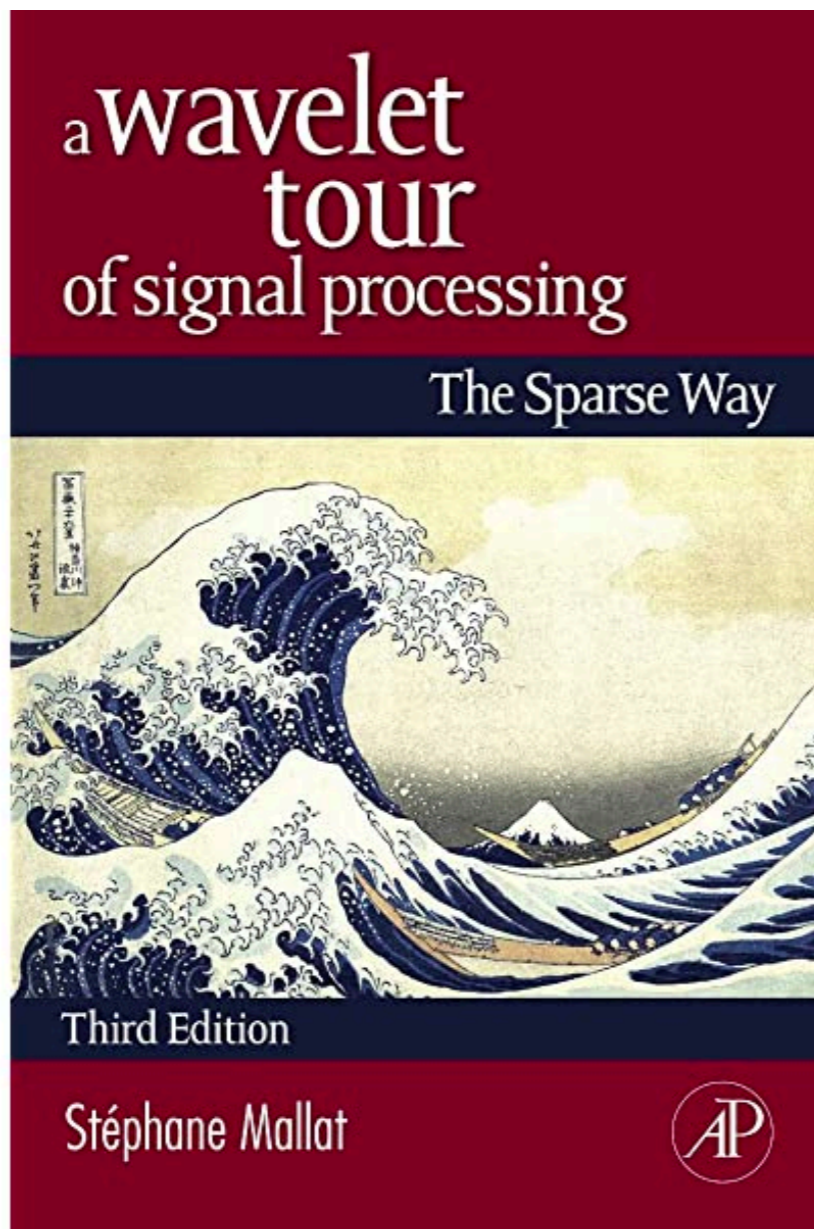
Sound Event Detection

Smart Home Technologies
Save Money and Lives

Harmonic Time Series

IEEE
Signal
Processing
Society





■ Some other references

Bach, F., Jenatton, R., Mairal, J., and Obozinski, G. Convex optimization with sparsity-inducing norms. *Foundations and Trends in Machine Learning*, 4(1):1–106, 2012.

Beck, A. and Teboulle, M. A fast iterative shrinkage thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.

Friedman, J., Hastie, T. J., Hofling, H., and Tibshirani, R. Pathwise coordinate optimization. *Ann. Appl. Stat.*, 1(2): 302–332, 2007.

E. J. Candes and T. Tao, Decoding by Linear Programming, *IEEE Trans. Inf. Th.*, 51(12): 4203–4215 (2005).

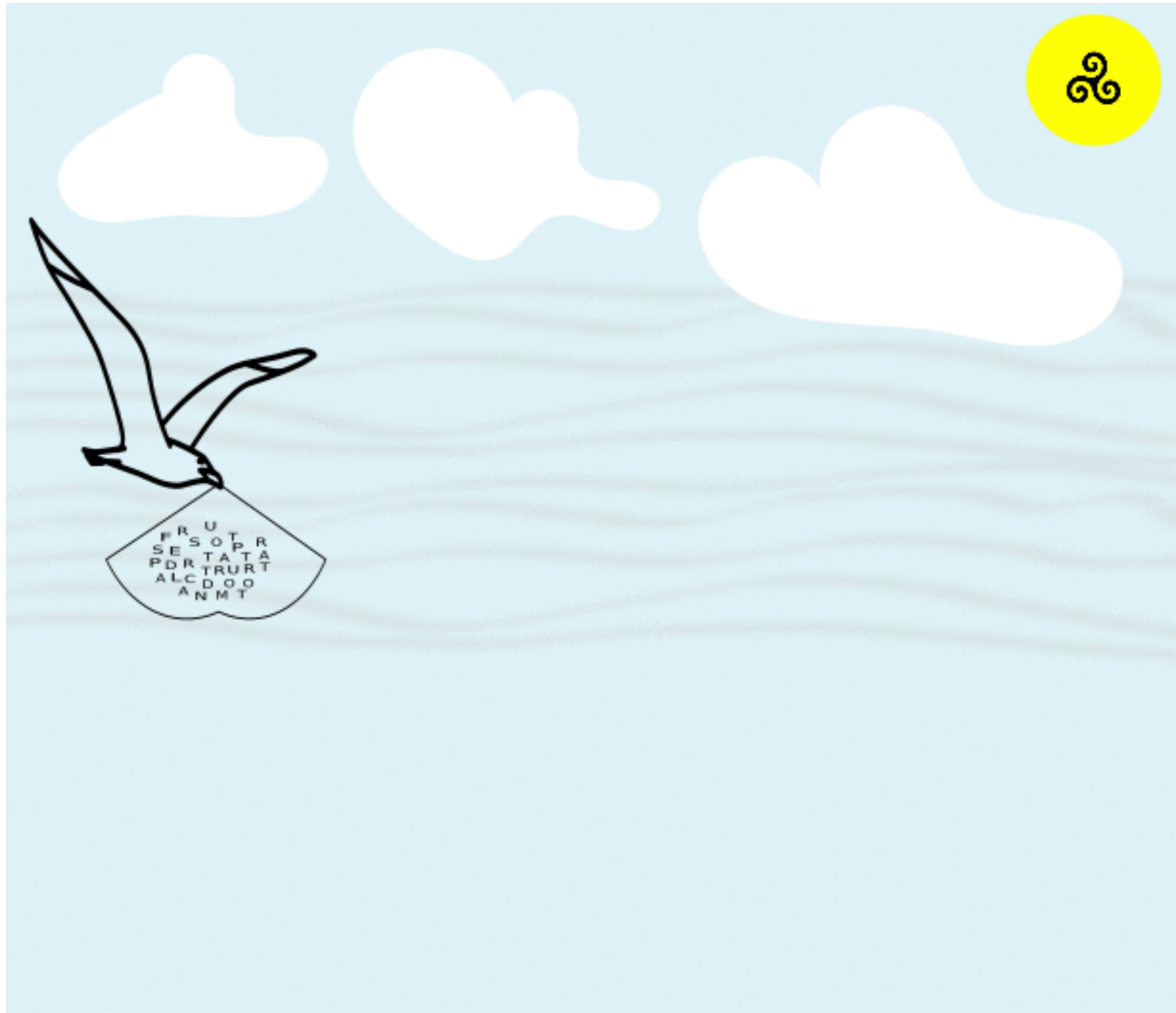
Arthur, D., Vassilvitskii, S. *k*-means++: the advantages of careful seeding . *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics Philadelphia, PA, USA. pp. 1027–1035.

Luc Giffon, Rémi Gribonval. Compressive Clustering with an Optical Processing Unit. GRETSI. 2022

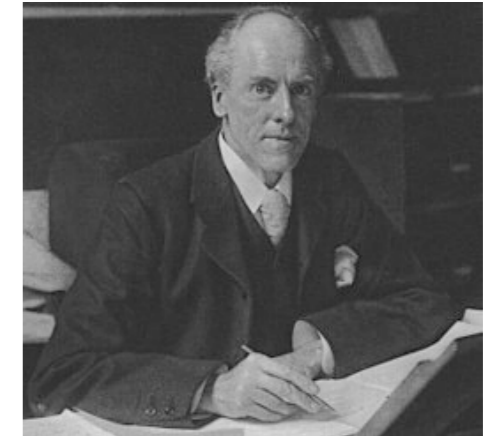
Chen, Scott Shaobing and Donoho, David L. and Saunders, Michael A. Atomic Decomposition by Basis Pursuit. *SIAM Journal on Scientific Computing*. 1995

Robert Tibshirani. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society*. 1996

Thank you!



Machine learning theory

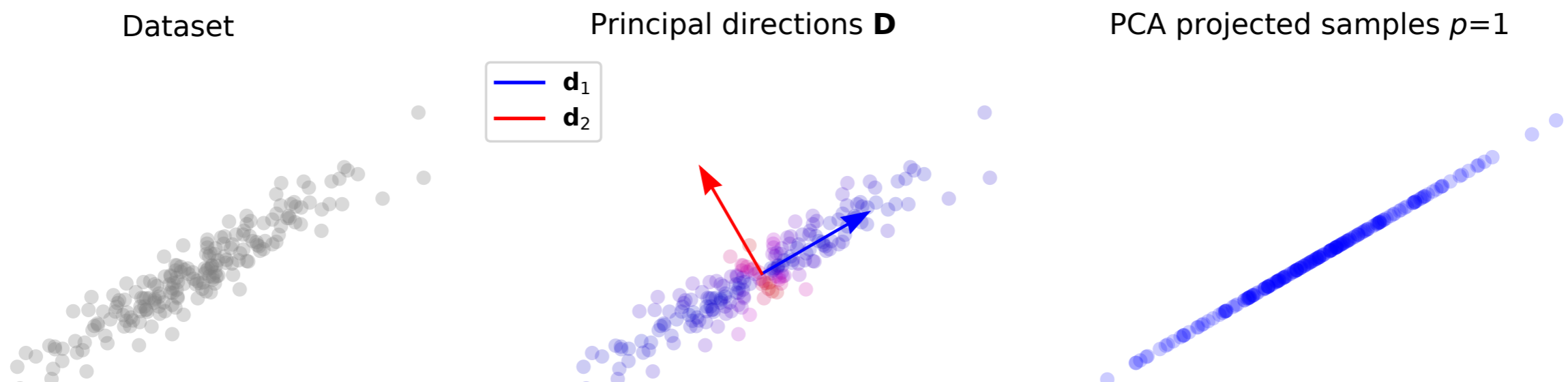


Dimension reduction: PCA

- We aim at solving: [Pearson, 1901]

$$\min_{\mathbf{D} \in \mathbb{R}^{d \times k}, \mathbf{D}^\top \mathbf{D} = \mathbf{I}_k} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\mathbf{D}^\top \mathbf{x}_i\|_2^2$$

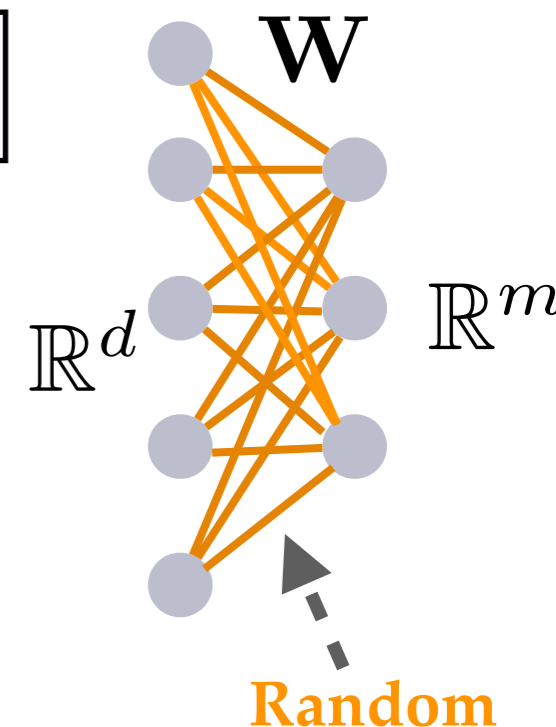
- Find a **linear subspace** that minimizes the **reconstruction error**
- $h = \mathbf{D}$, $\mathcal{H} = \text{Stiefel}$ and $\ell(\mathbf{x}_i, h) = \|\mathbf{x}_i - \mathbf{D}\mathbf{D}^\top \mathbf{x}_i\|_2^2$
- Eigen. decomposition of covariance matrix: $\mathcal{O}(nd^2 + d^3)$



Theory of sketching

Randomization: the core of sketching

- A function called **feature operator** $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$
- In practice $\Phi(\mathbf{x}) = \rho(\mathbf{W}\mathbf{x})$
- Where $\mathbf{W} \in \mathbb{R}^{m \times d}$ is a **random matrix**
- Where $\rho : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is a **non-linear activation**



Random Fourier Features (RFF)

- Where $\mathbf{W} \in \mathbb{R}^{m \times d}$ is **Gaussian** $W_{ij} \sim \mathcal{N}(0, \sigma^2)$
- Where $\rho(\mathbf{y}) = \frac{1}{\sqrt{m}} (\exp(-iy_1), \dots, \exp(-iy_m))$

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{m}} (\exp(-i\boldsymbol{\omega}_1^\top \mathbf{x}), \dots, \exp(-i\boldsymbol{\omega}_m^\top \mathbf{x}))$$

$$\mathbf{W} = [\boldsymbol{\omega}_1^\top, \dots, \boldsymbol{\omega}_m^\top]$$

- For PCA:

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{m}} (|\mathbf{x}^\top \boldsymbol{\omega}_1|^2, \dots, |\mathbf{x}^\top \boldsymbol{\omega}_m|^2)$$

[Rahimi & Recht, 2008]

