

Random Moments for Sketched Mixture Learning

Nicolas Keriven^{1,2}, Rémi Gribonval²,
Gilles Blanchard³, Yann Traonmilin²

¹Université Rennes 1

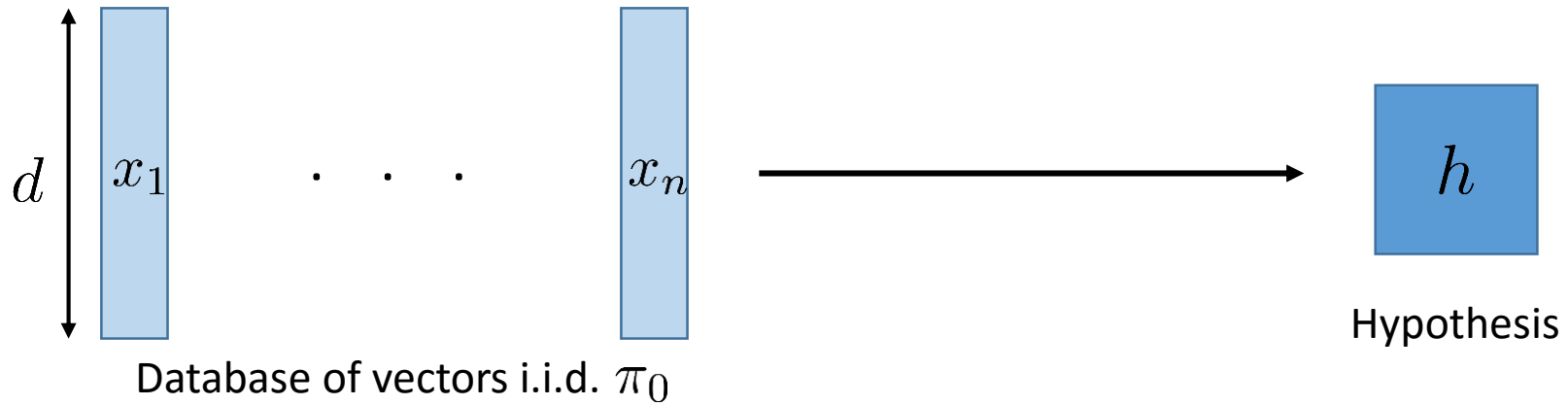
²Inria Rennes Bretagne-atlantique

³University of Potsdam

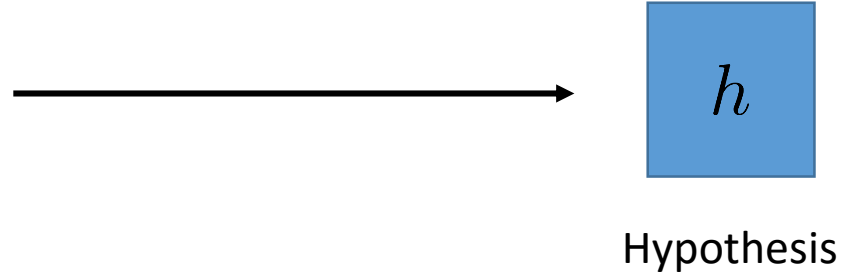
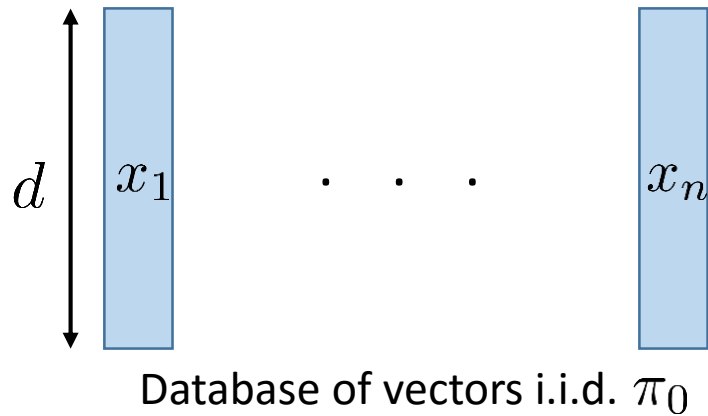
SPARS 2017

- 1** Introduction
- 2 Illustration
- 3 Main results
- 4 Conclusion

Statistical Learning

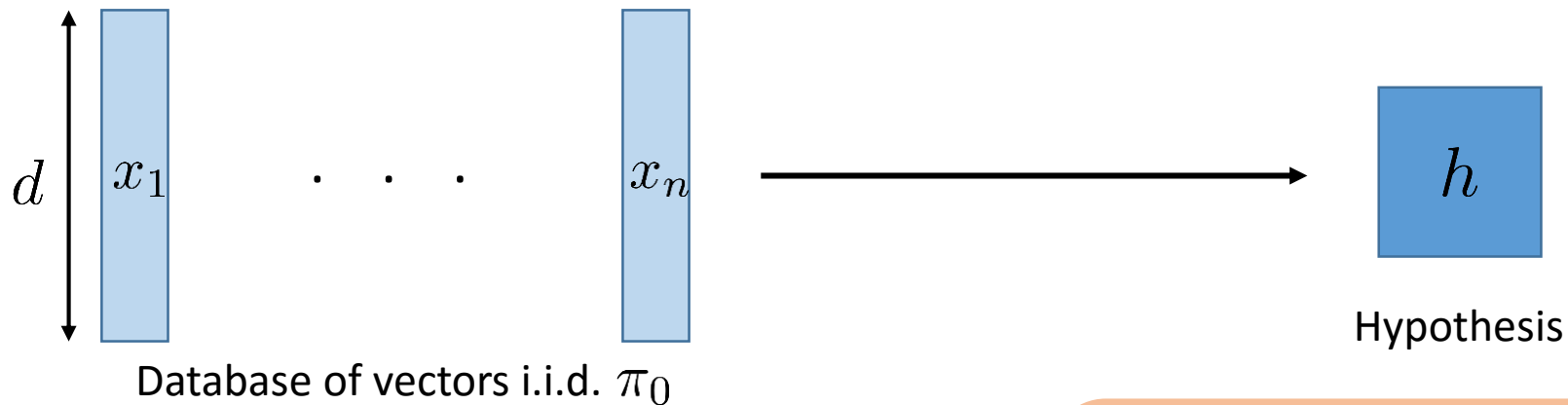


Statistical Learning



- PCA : $\mathbf{x} \in \text{Span}(h_1, \dots, h_k)$
- Classification : $\langle h, \Phi(\mathbf{x}) \rangle$
- Regression : $\mathbf{y} = h(\mathbf{x})$
- **k-means** : $h = \{c_l\}_{l=1}^k$
- **Density estimation** : $\mathbf{x} \sim \pi_h$

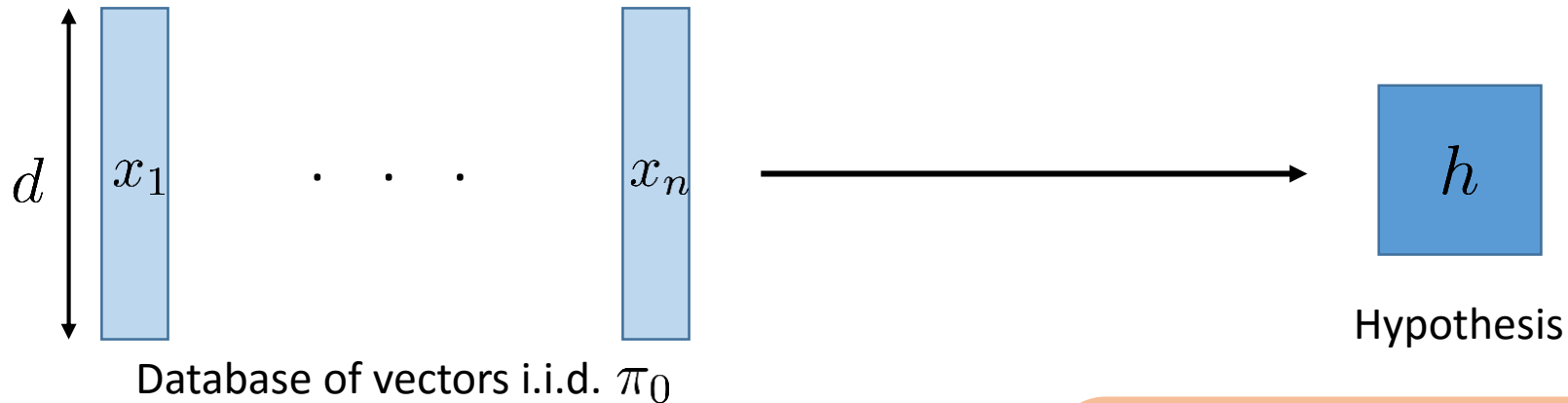
Statistical Learning



Loss function $\ell(x, h) \in \mathbb{R}$

- PCA : $\mathbf{x} \in \text{Span}(h_1, \dots, h_k)$
- Classification : $\langle h, \Phi(\mathbf{x}) \rangle$
- Regression : $\mathbf{y} = h(\mathbf{x})$
- **k-means** : $h = \{c_l\}_{l=1}^k$
- **Density estimation** : $\mathbf{x} \sim \pi_h$

Statistical Learning



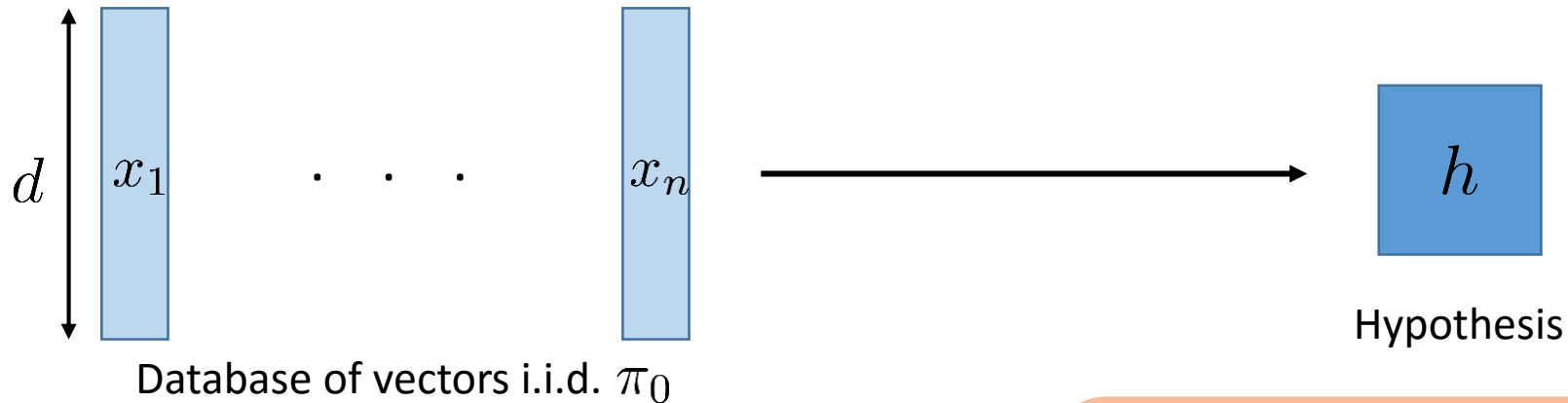
Loss function $\ell(x, h) \in \mathbb{R}$

Goal : Minimize Expected Risk

$$h^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

- PCA : $\mathbf{x} \in \text{Span}(h_1, \dots, h_k)$
- Classification : $\langle h, \Phi(\mathbf{x}) \rangle$
- Regression : $\mathbf{y} = h(\mathbf{x})$
- k-means : $h = \{c_l\}_{l=1}^k$
- Density estimation : $\mathbf{x} \sim \pi_h$

Statistical Learning



Loss function $\ell(x, h) \in \mathbb{R}$

Goal : Minimize Expected Risk

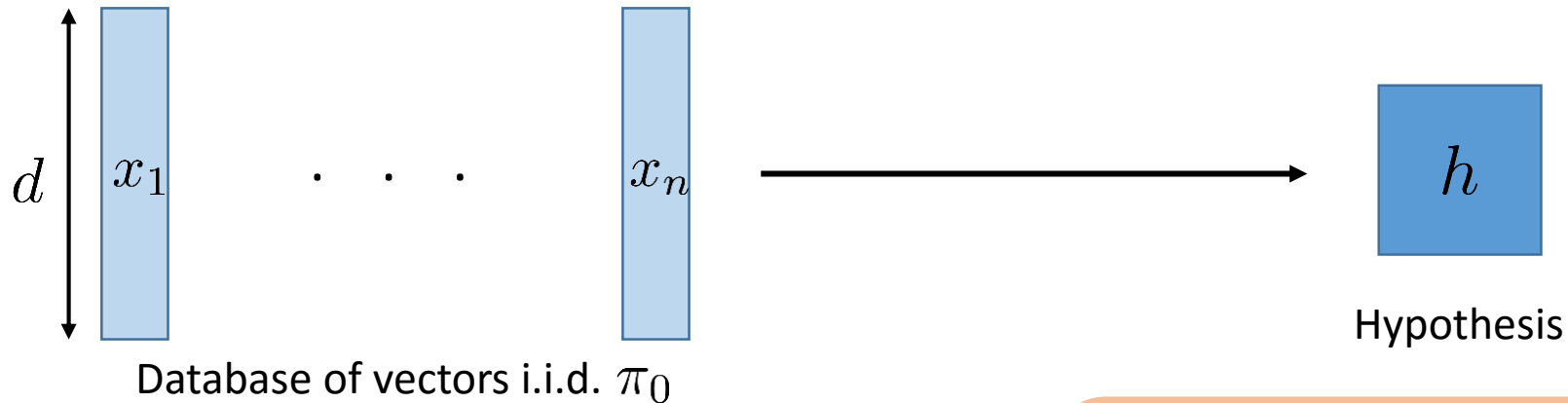
$$h^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

Empirical Risk Minimization (ERM)

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(x_i, h)$$

- PCA : $\mathbf{x} \in \text{Span}(h_1, \dots, h_k)$
- Classification : $\langle h, \Phi(\mathbf{x}) \rangle$
- Regression : $\mathbf{y} = h(\mathbf{x})$
- k-means : $h = \{c_l\}_{l=1}^k$
- Density estimation : $\mathbf{x} \sim \pi_h$

Statistical Learning



Loss function $\ell(x, h) \in \mathbb{R}$

Goal : Minimize Expected Risk

$$h^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

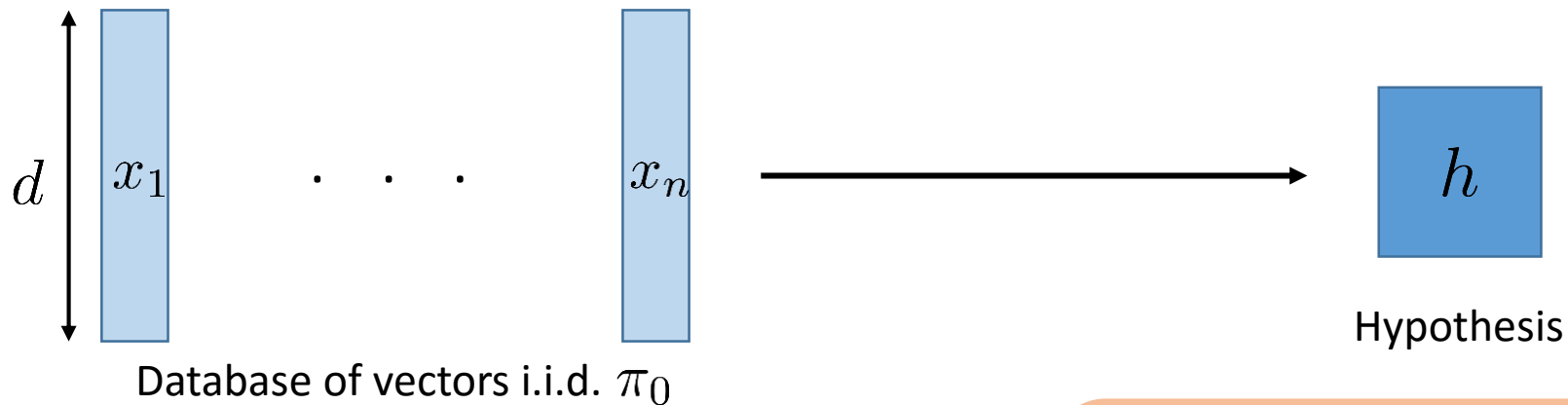
Empirical Risk Minimization (ERM)

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(x_i, h)$$

- PCA : $\mathbf{x} \in \text{Span}(h_1, \dots, h_k)$
- Classification : $\langle h, \Phi(\mathbf{x}) \rangle$
- Regression : $\mathbf{y} = h(\mathbf{x})$
- k-means : $h = \{c_l\}_{l=1}^k$
- Density estimation : $\mathbf{x} \sim \pi_h$

Large d or n

Statistical Learning



Loss function $\ell(x, h) \in \mathbb{R}$

Goal : Minimize Expected Risk

$$h^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

Empirical Risk Minimization (ERM)

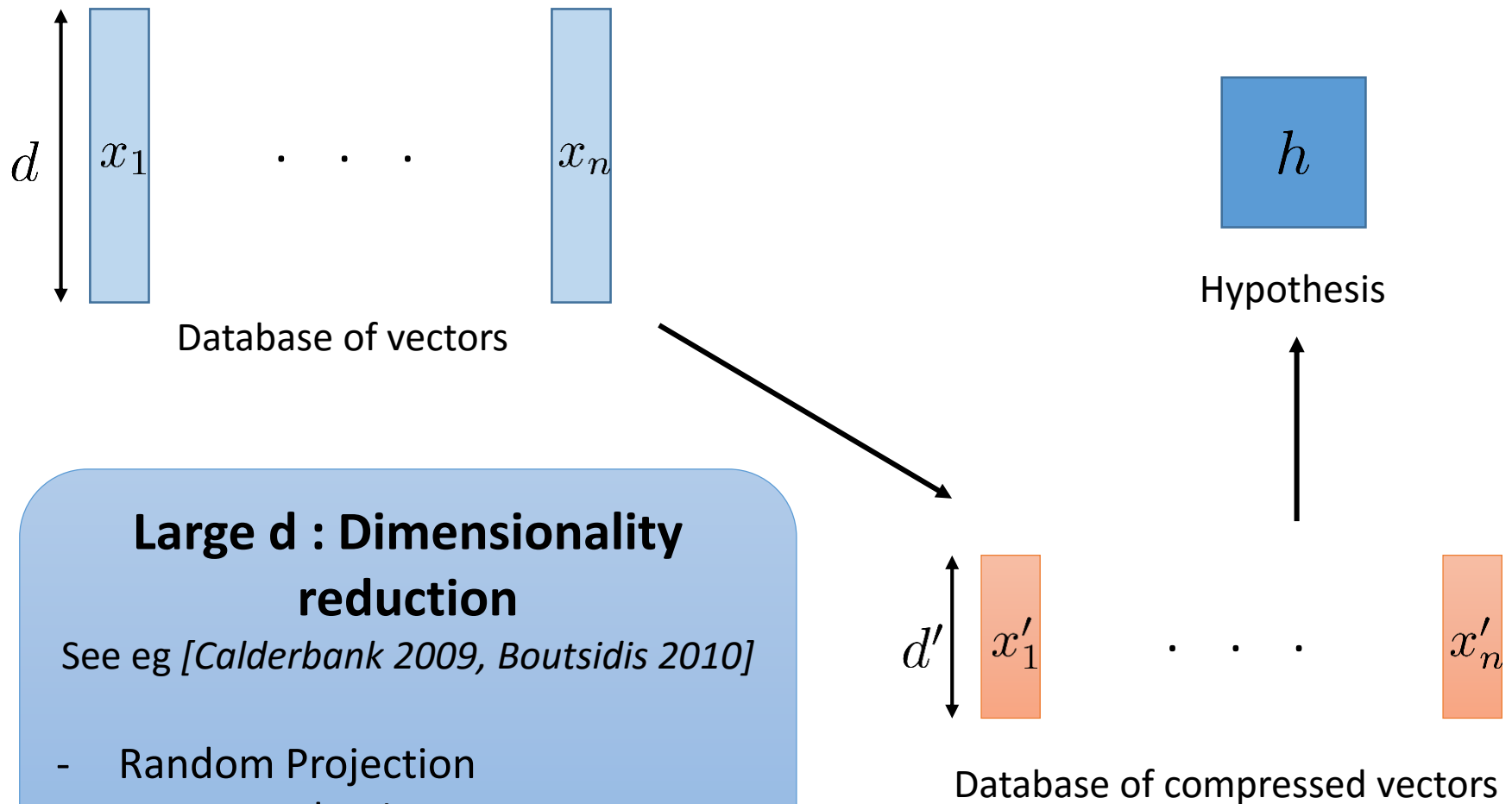
$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(x_i, h)$$

- PCA : $\mathbf{x} \in \text{Span}(h_1, \dots, h_k)$
- Classification : $\langle h, \Phi(\mathbf{x}) \rangle$
- Regression : $\mathbf{y} = h(\mathbf{x})$
- k-means : $h = \{c_l\}_{l=1}^k$
- Density estimation : $\mathbf{x} \sim \pi_h$

Large d or n

Compress the database before learning

Compressive Statistical Learning

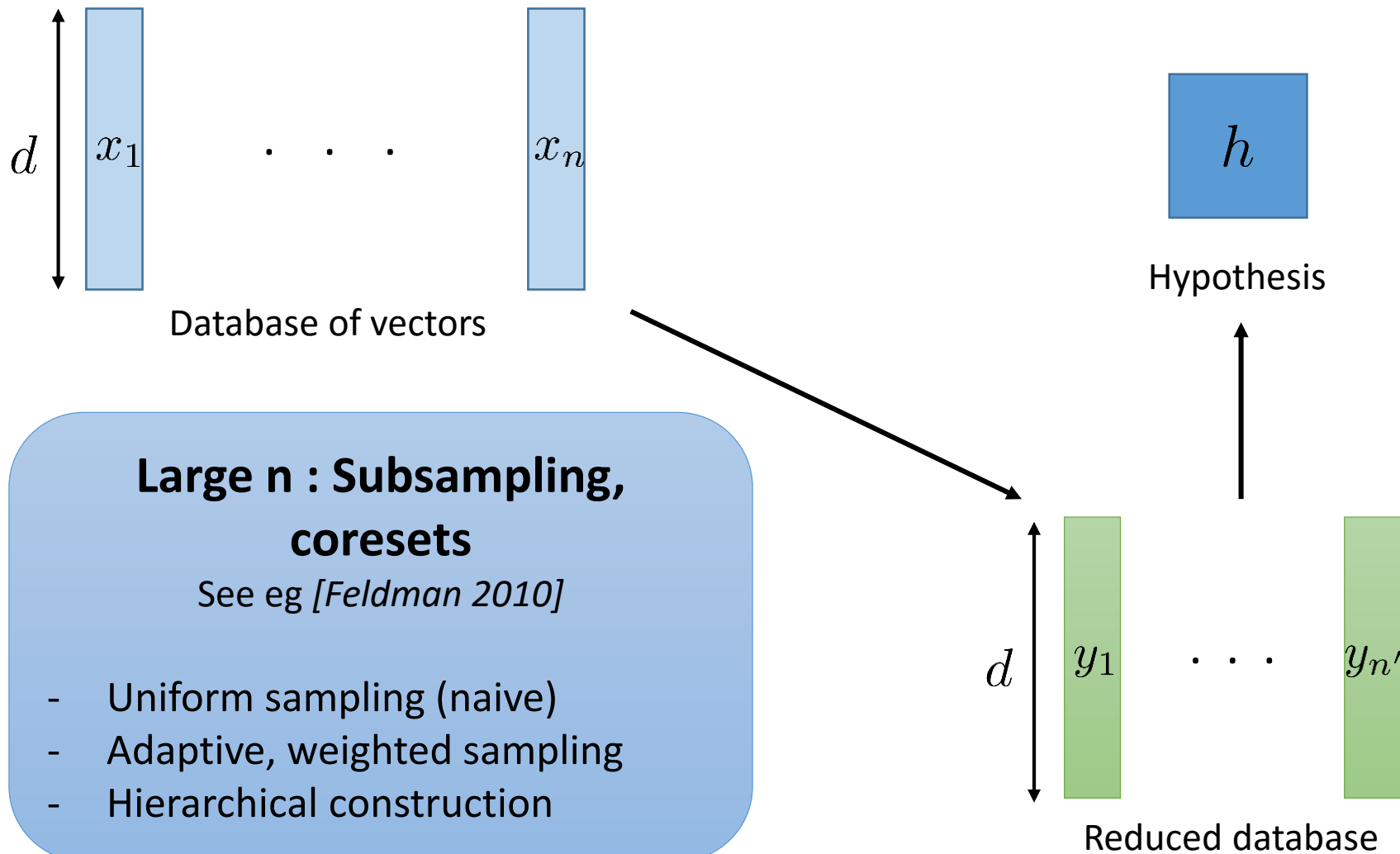


Large d : Dimensionality reduction

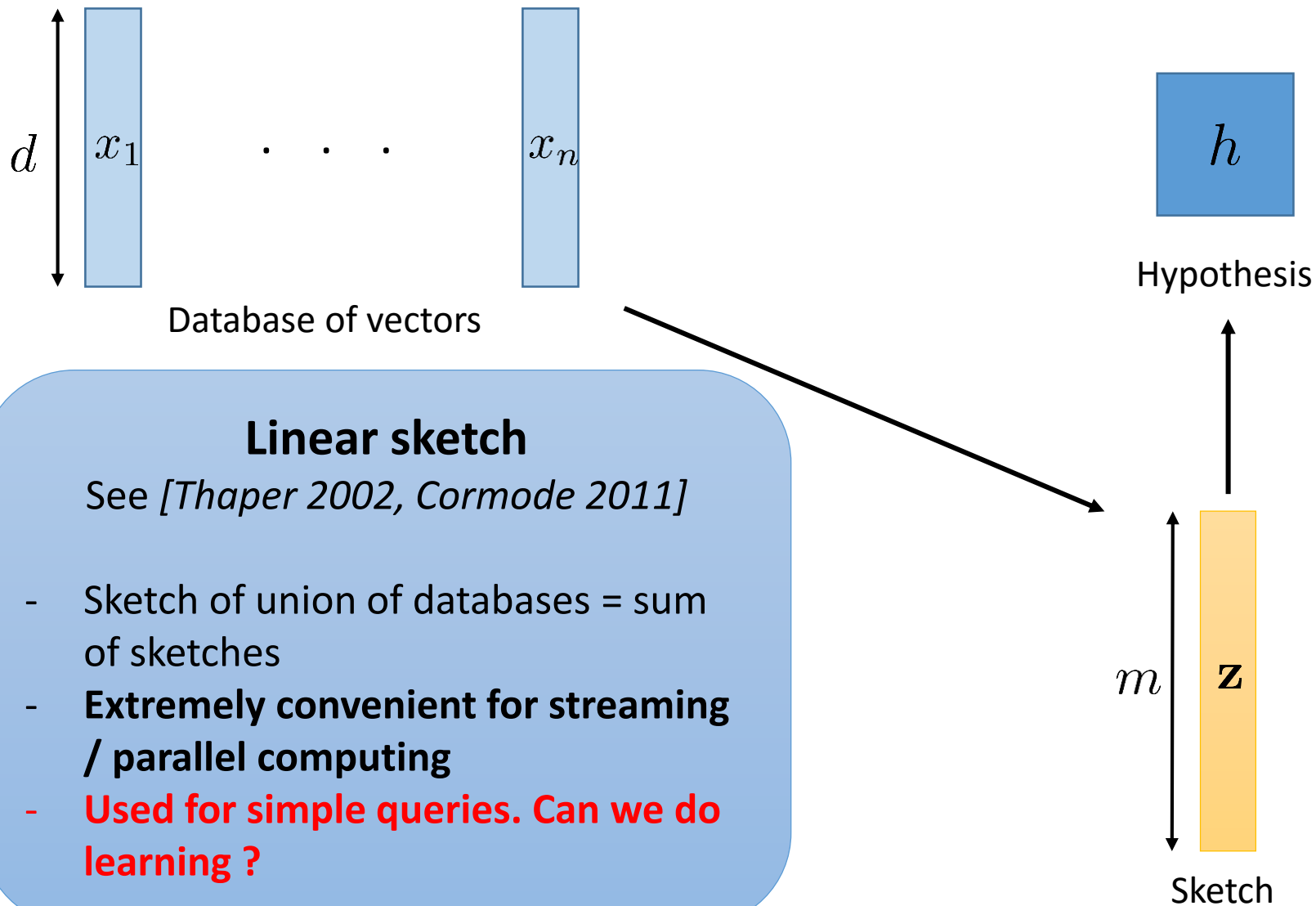
See eg [Calderbank 2009, Boutsidis 2010]

- Random Projection
- Feature selection

Compressive Statistical Learning



Compressive Statistical Learning



Random Sketching operator

**Linear sketch = Empirical
generalized moments...**

$$\mathbf{z} = \left(\frac{1}{n}\right) \sum_{i=1}^n \Phi(x_i)$$

Random Sketching operator

Linear sketch = Empirical
generalized moments...

$$\mathbf{z} = \left(\frac{1}{n}\right) \sum_{i=1}^n \Phi(x_i)$$

... i.e. *linear measurements of
underlying probability
distribution*

$$\mathbf{z} \approx \mathbb{E}_{x \sim \pi_0} \Phi(x) = \mathcal{A}\pi_0$$

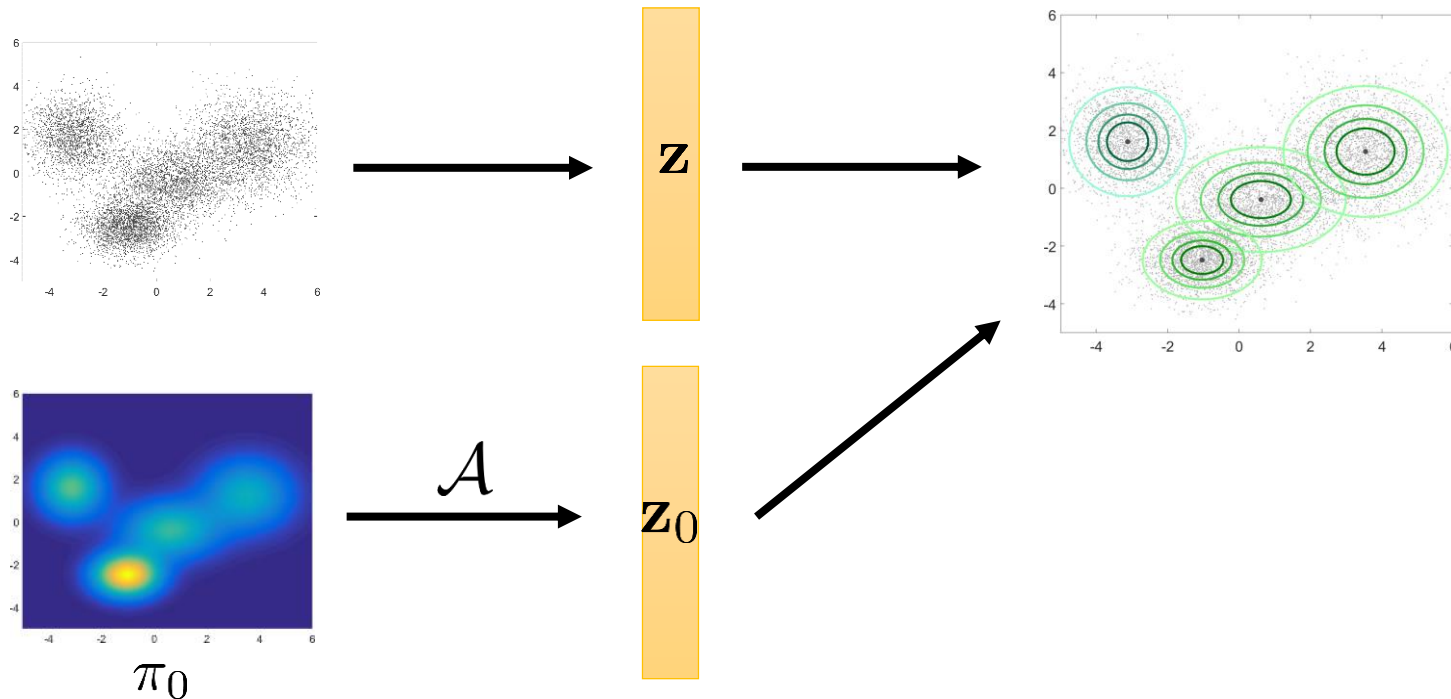
Random Sketching operator

Linear sketch = Empirical generalized moments...

$$\mathbf{z} = \left(\frac{1}{n}\right) \sum_{i=1}^n \Phi(x_i)$$

... i.e. *linear measurements of underlying probability distribution*

$$\mathbf{z} \approx \mathbb{E}_{x \sim \pi_0} \Phi(x) = \mathcal{A}\pi_0$$



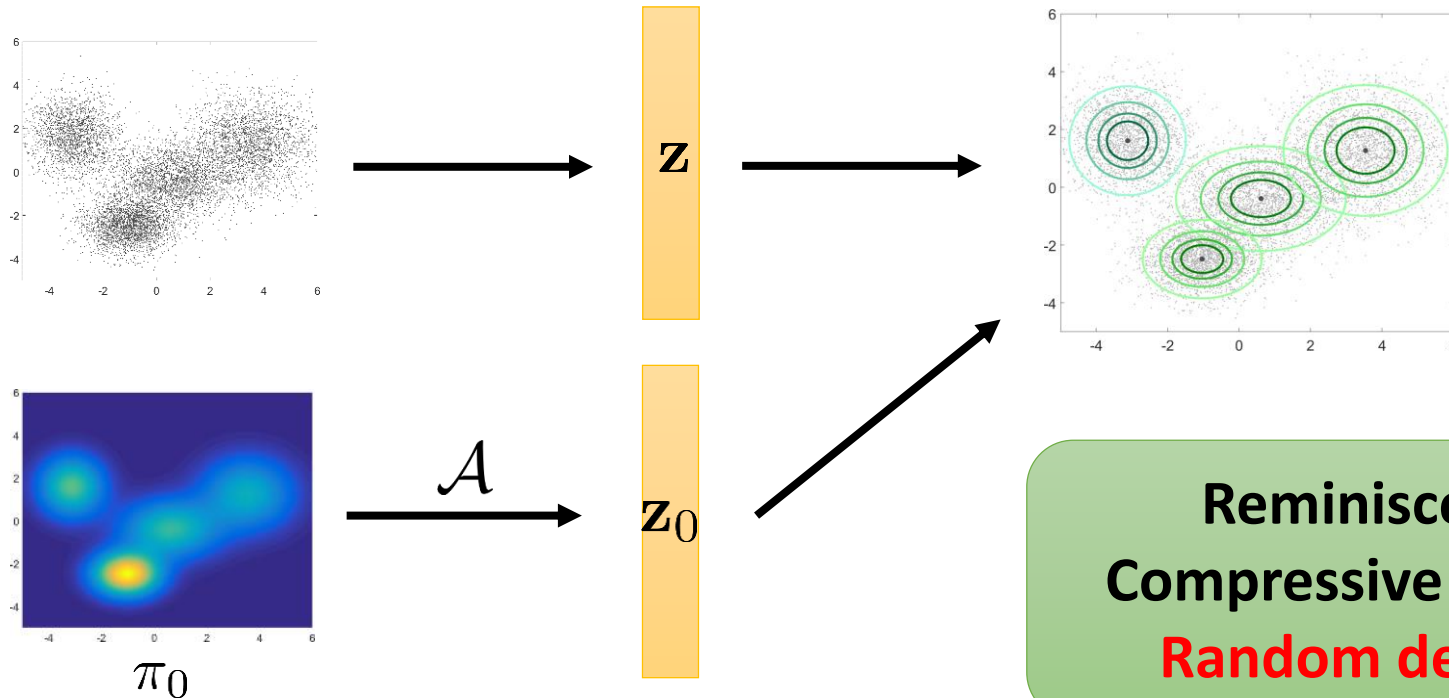
Random Sketching operator

Linear sketch = Empirical generalized moments...

$$\mathbf{z} = \left(\frac{1}{n}\right) \sum_{i=1}^n \Phi(x_i)$$

... i.e. *linear measurements of underlying probability distribution*

$$\mathbf{z} \approx \mathbb{E}_{x \sim \pi_0} \Phi(x) = \mathcal{A}\pi_0$$



Reminiscent of
Compressive Sensing :
Random design of \mathcal{A}

Outline

- ① Introduction
- ② **Illustration (previous work)**
- ③ Main results
- ④ Conclusion

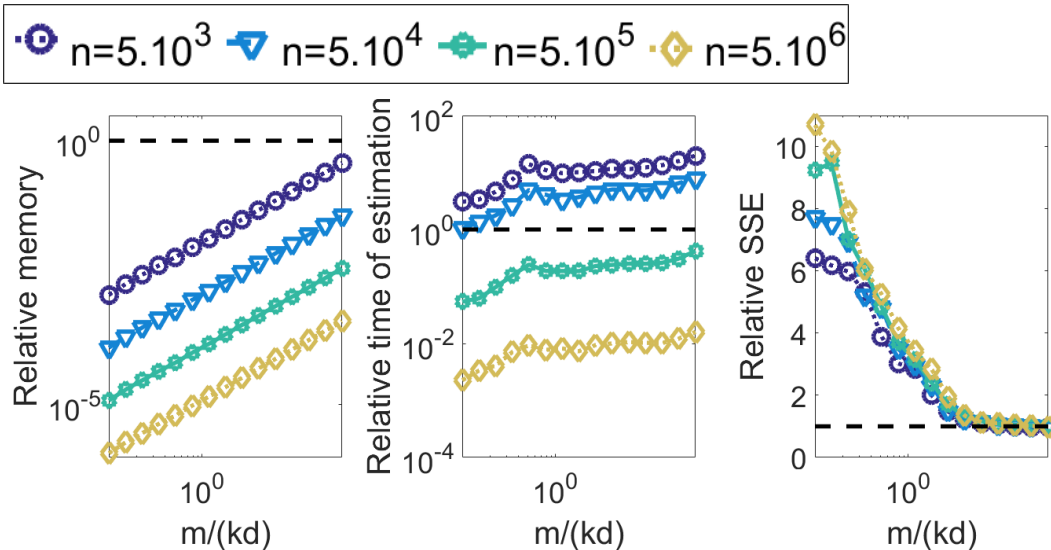
Experimental illustration

Compressive Learning-OMP algorithm [*Keriven 2015,2016*]
(OMP + non-convex updates)

Experimental illustration

k-means (d=10, k=10)

Compressive Learning-OMP algorithm [Keriven 2015,2016]
(OMP + non-convex updates)



Comparison with

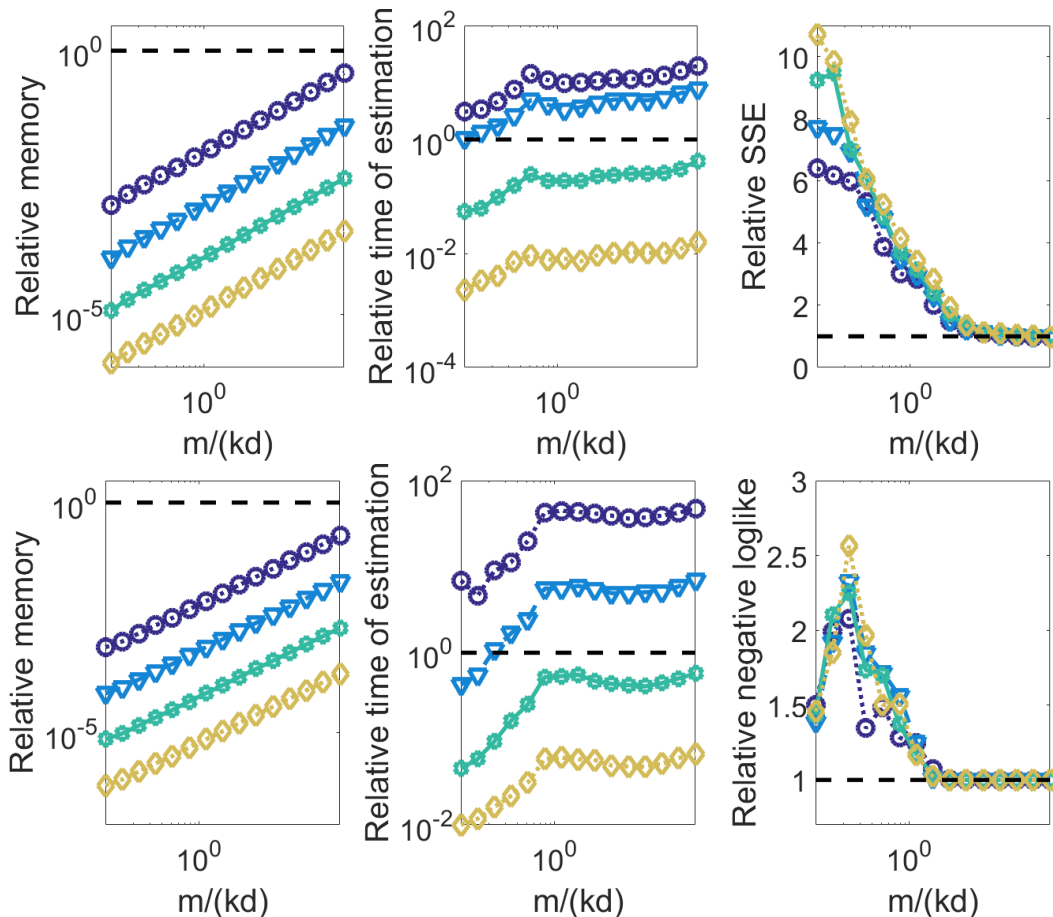
- *Matlab's* kmeans
- *VLFeat's* gmm
- Faster and more memory efficient on large databases
- Number of measurements does not depend on n

Experimental illustration

k-means (d=10, k=10)

Compressive Learning-OMP algorithm [Keriven 2015,2016]
(OMP + non-convex updates)

• $n=5 \cdot 10^3$ • $n=5 \cdot 10^4$ • $n=5 \cdot 10^5$ • $n=5 \cdot 10^6$



Comparison with

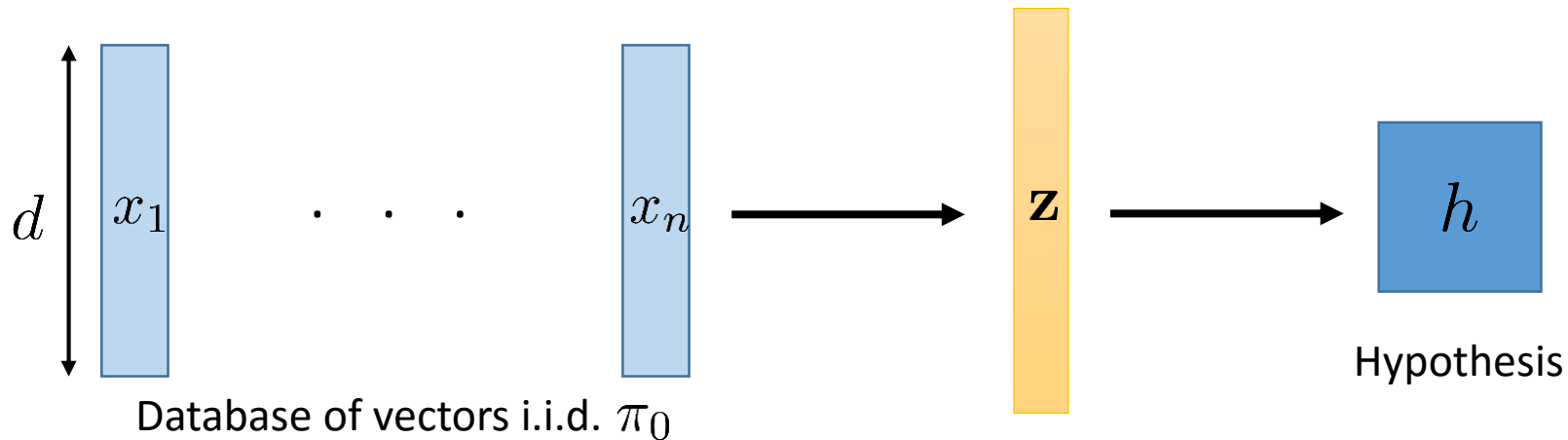
- *Matlab's* kmeans
- *VLFeat's* gmm
- Faster and more memory efficient on large databases
- Number of measurements does not depend on n

GMMs (d=10, k=10)

Outline

- ① Introduction
- ② Illustration
- ③ **Main results**
- ④ Conclusion

Statistical Learning



Loss function

$$\ell(x, h)$$

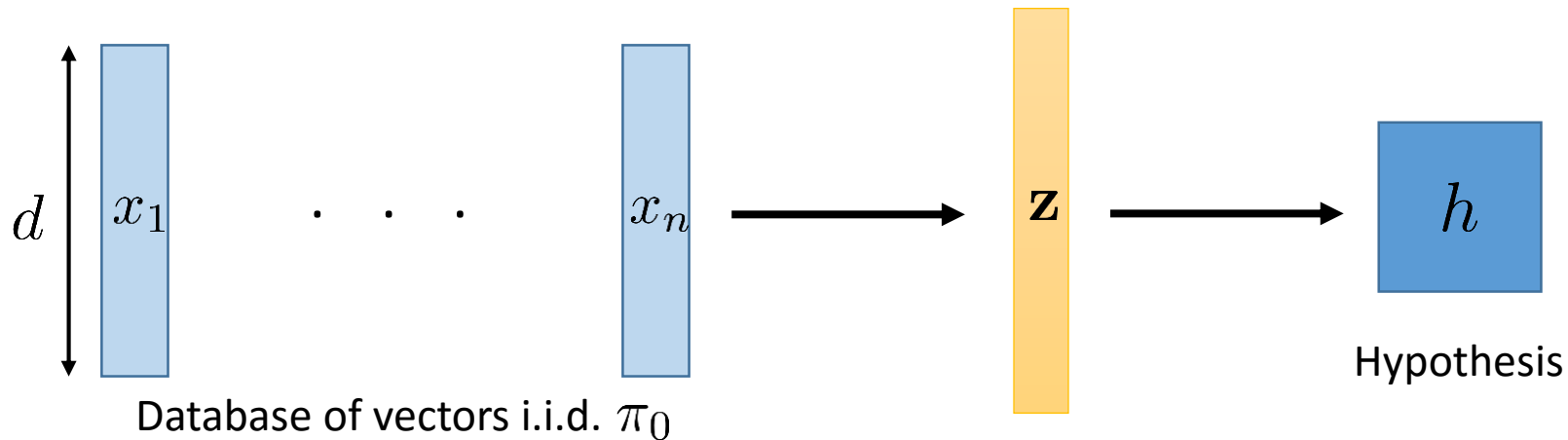
Goal : Minimize Expected Risk

$$h^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

$$\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \Phi(x_i)$$

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

Statistical Learning



Loss function

$$\ell(x, h)$$

Goal : Minimize Expected Risk

$$h^* = \arg \min_{h \in \mathcal{H}} \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

$$\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \Phi(x_i)$$

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

Here:

- k-means
- GMM with known covariance

Hyp. class $h = \{c_1, \dots, c_k\} \subset \mathbb{R}^d$

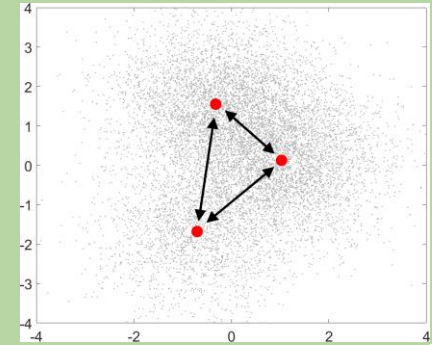
k-means

Hyp. class $h = \{c_1, \dots, c_k\} \subset \mathbb{R}^d$

$\mathcal{H}_{k, \varepsilon, M}$

- M - bounded domain
(centroids, not samples)

- ε - separation



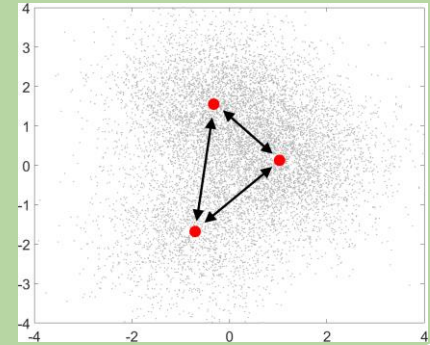
k-means

Hyp. class $h = \{c_1, \dots, c_k\} \subset \mathbb{R}^d$

$\mathcal{H}_{k, \varepsilon, M}$

- M - bounded domain
(centroids, not samples)

- ε - separation



Loss function $\ell(x, h) = \min_{1 \leq l \leq k} \|x - c_l\|_2^2$

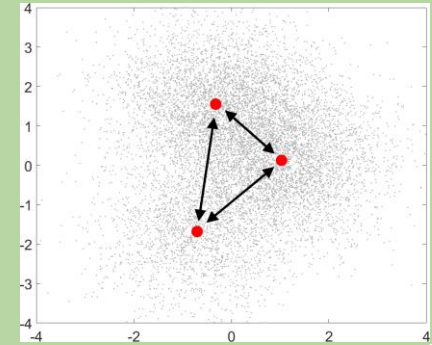
k-means

Hyp. class $h = \{c_1, \dots, c_k\} \subset \mathbb{R}^d$

- ε - separation

$\mathcal{H}_{k,\varepsilon,M}$

- M - bounded domain
(centroids, not samples)



Loss function $\ell(x, h) = \min_{1 \leq l \leq k} \|x - c_l\|_2^2$

Sketching operator

- (weighted) Random Fourier sampling
- « Smoothing » weights $c_\omega \propto \|\omega\|_2^2$

$\{\omega_1, \dots, \omega_m\} \subset (\mathbb{R}^d)^m$

$$\Phi(x) = \left[e^{-i\omega_j^T x} / c_{\omega_j} \right]_{j=1}^m$$

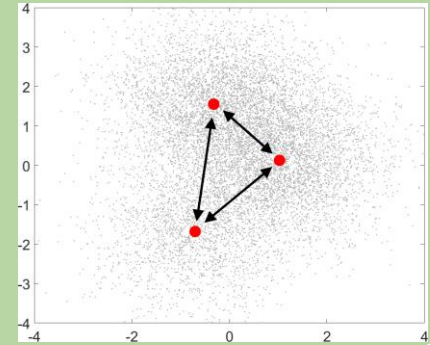
k-means

Hyp. class $h = \{c_1, \dots, c_k\} \subset \mathbb{R}^d$

- ε - separation

$\mathcal{H}_{k,\varepsilon,M}$

- M - bounded domain
(centroids, not samples)



Loss function $\ell(x, h) = \min_{1 \leq l \leq k} \|x - c_l\|_2^2$

Sketching operator

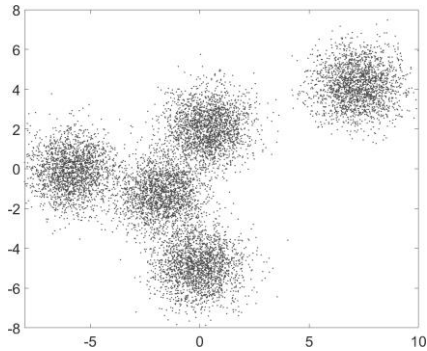
$\{\omega_1, \dots, \omega_m\} \subset (\mathbb{R}^d)^m$

$$\Phi(x) = \left[e^{-i\omega_j^T x} / c_{\omega_j} \right]_{j=1}^m$$

- (weighted) Random Fourier sampling
- « Smoothing » weights $c_\omega \propto \|\omega\|_2^2$

$$\omega_j \stackrel{i.i.d.}{\sim} \Lambda(\omega) \propto c_\omega^2 \mathcal{N}(0, \sigma^2 \mathbf{Id})$$
$$\sigma^2 \propto \varepsilon^{-1}$$

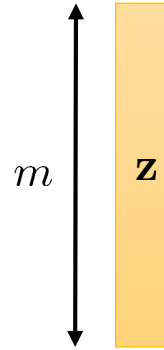
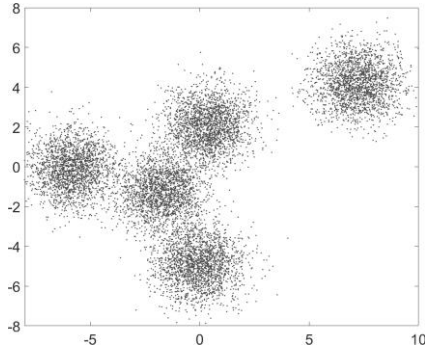
k-means: result



$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

k-means: result



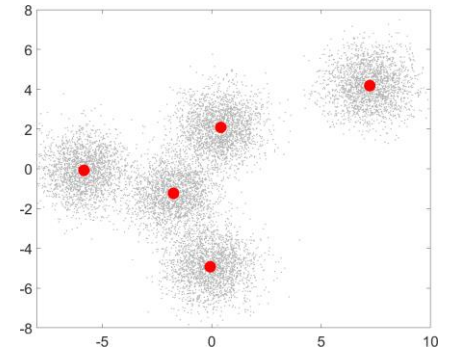
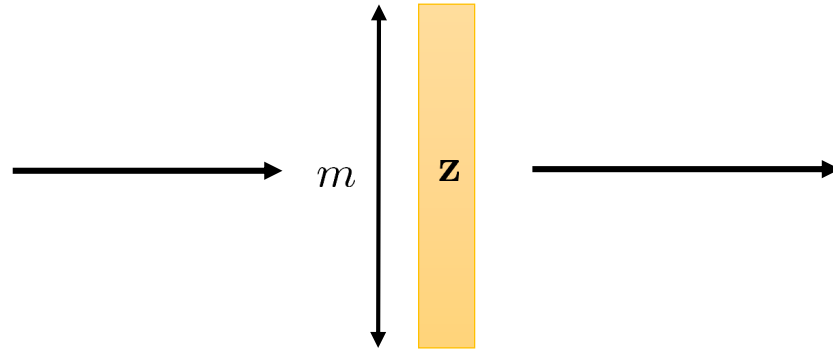
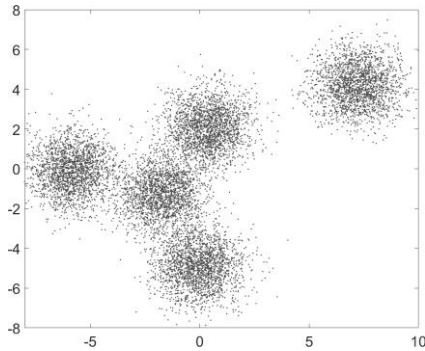
$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

$$\omega_1, \dots, \omega_m \stackrel{i.i.d.}{\sim} \Lambda$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

Z: (weighted) Random Fourier sampling

k-means: result



$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

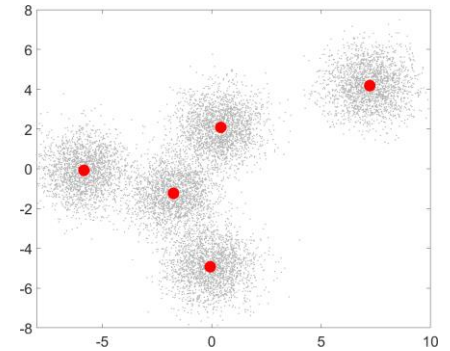
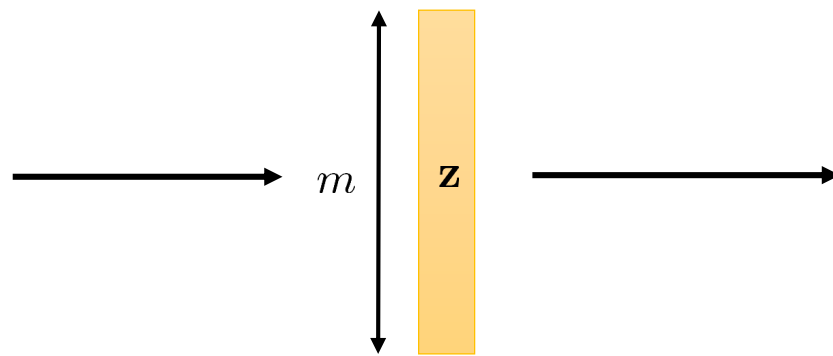
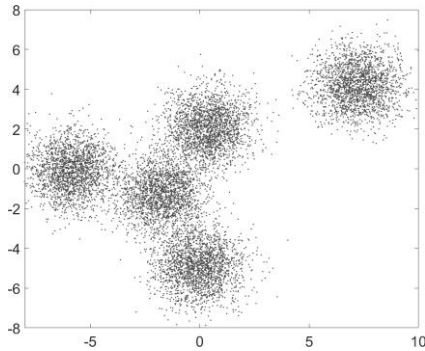
$$\omega_1, \dots, \omega_m \stackrel{i.i.d.}{\sim} \Lambda$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

\mathbf{Z} : (weighted) Random Fourier sampling

$$\hat{h} = \min_{h \in \mathcal{H}_{k, \varepsilon, M}} \min_{\alpha \geq 0, \sum_l \alpha_l = 1} \|\mathbf{z} - \mathcal{A}(\sum_{l=1}^k \alpha_l \delta_{c_l})\|_2$$

k-means: result



$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

$$\omega_1, \dots, \omega_m \stackrel{i.i.d.}{\sim} \Lambda$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

z: (weighted) Random Fourier sampling

$$\hat{h} = \min_{h \in \mathcal{H}_{k, \varepsilon, M}} \min_{\alpha \geq 0, \sum_l \alpha_l = 1} \|\mathbf{z} - \mathcal{A}(\sum_{l=1}^k \alpha_l \delta_{c_l})\|_2$$

$$h^* = \min_{h \in \mathcal{H}_{k, \varepsilon, M}} R_{\pi_0}(h)$$

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

w.h.p. on x_i, ω_j

$$R_{\pi_0}(\hat{h}) \lesssim R_{\pi_0}(h^*) + \mathcal{O}(\sqrt{1/n})$$

GMM with known covariance Σ

Hyp. class $h = \{(\mu_1, \alpha_1), \dots, (\mu_k, \alpha_k)\} \subset \mathbb{R}^d \times \mathbb{R}_+$

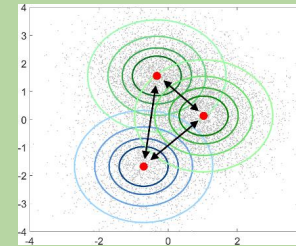
GMM with known covariance Σ

Hyp. class $h = \{(\mu_1, \alpha_1), \dots, (\mu_k, \alpha_k)\} \subset \mathbb{R}^d \times \mathbb{R}_+$

- $\varepsilon \geq \varepsilon_0$ separation

$\mathcal{H}_{k, \varepsilon, M}$

- M - bounded domain
(means, not samples)



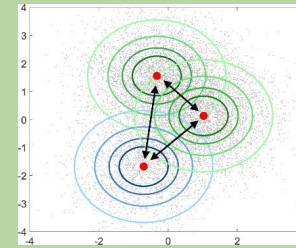
GMM with known covariance Σ

Hyp. class $h = \{(\mu_1, \alpha_1), \dots, (\mu_k, \alpha_k)\} \subset \mathbb{R}^d \times \mathbb{R}_+$

- $\varepsilon \geq \varepsilon_0$ separation

$\mathcal{H}_{k, \varepsilon, M}$

- M - bounded domain
(means, not samples)



Loss function $\pi_h = \sum_{l=1}^k \alpha_l \mathcal{N}(\mu_l, \Sigma)$ $\ell(x, h) = -\log \pi_h(x)$

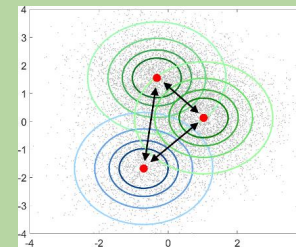
GMM with known covariance Σ

Hyp. class $h = \{(\mu_1, \alpha_1), \dots, (\mu_k, \alpha_k)\} \subset \mathbb{R}^d \times \mathbb{R}_+$

- $\varepsilon \geq \varepsilon_0$ separation

$\mathcal{H}_{k, \varepsilon, M}$

- M - bounded domain
(means, not samples)



Loss function $\pi_h = \sum_{l=1}^k \alpha_l \mathcal{N}(\mu_l, \Sigma)$ $\ell(x, h) = -\log \pi_h(x)$

Sketching operator

- Random Fourier sampling

$$\{\omega_1, \dots, \omega_m\} \subset (\mathbb{R}^d)^m$$

$$\Phi(x) = \left[e^{-i\omega_j^T x} \right]_{j=1}^m$$

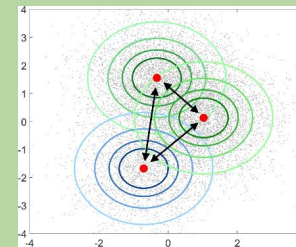
GMM with known covariance Σ

Hyp. class $h = \{(\mu_1, \alpha_1), \dots, (\mu_k, \alpha_k)\} \subset \mathbb{R}^d \times \mathbb{R}_+$

- $\varepsilon \geq \varepsilon_0$ separation

$\mathcal{H}_{k, \varepsilon, M}$

- M - bounded domain
(means, not samples)



Loss function $\pi_h = \sum_{l=1}^k \alpha_l \mathcal{N}(\mu_l, \Sigma)$ $\ell(x, h) = -\log \pi_h(x)$

Sketching operator

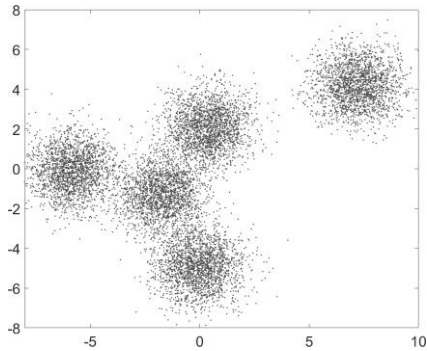
- Random Fourier sampling

$$\{\omega_1, \dots, \omega_m\} \subset (\mathbb{R}^d)^m$$

$$\Phi(x) = \left[e^{-i\omega_j^T x} \right]_{j=1}^m$$

σ^2 linked to separation
 $\omega_j \stackrel{i.i.d.}{\sim} \Lambda(\omega) = \mathcal{N}(0, \sigma^2 \Sigma^{-1})$

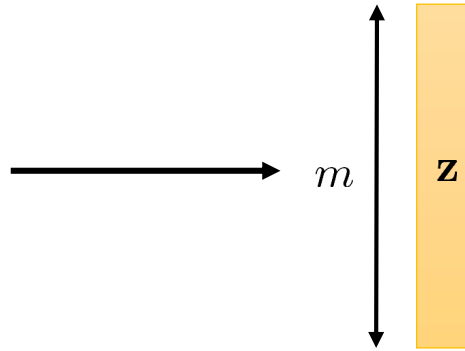
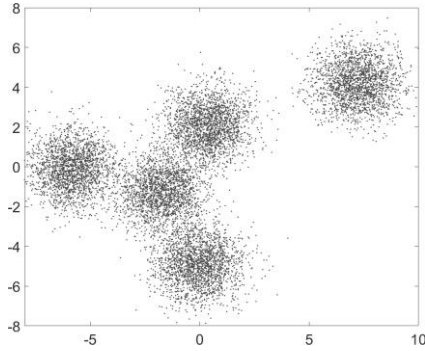
GMM: result



$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

GMM: result



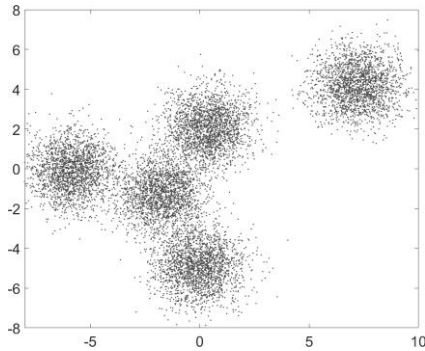
$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

$$\omega_1, \dots, \omega_m \stackrel{i.i.d.}{\sim} \Lambda$$

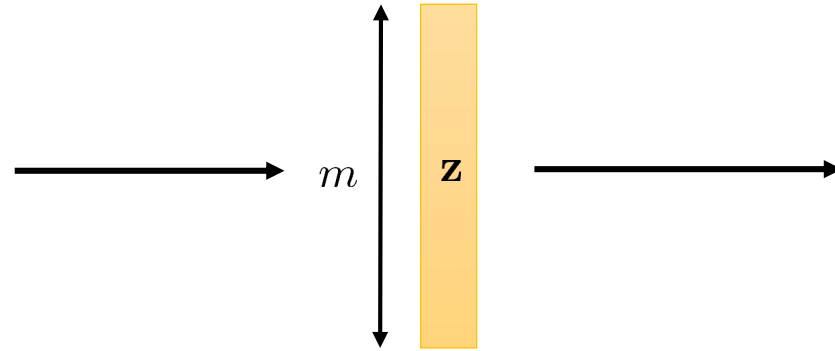
Z: Random Fourier sampling

GMM: result



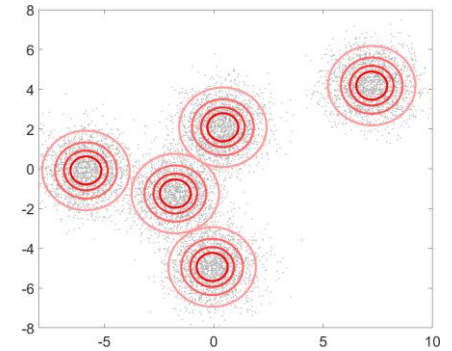
$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$



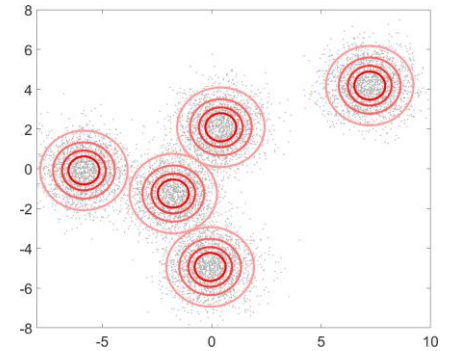
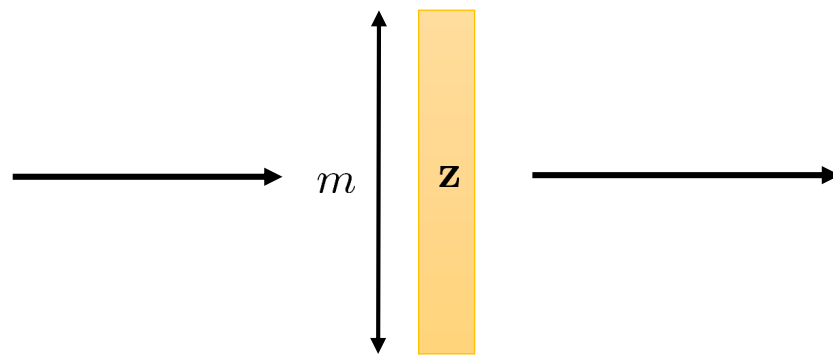
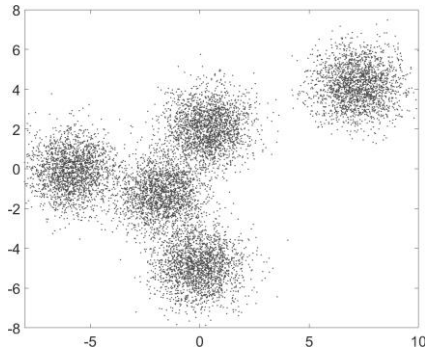
$$\omega_1, \dots, \omega_m \stackrel{i.i.d.}{\sim} \Lambda$$

z: Random Fourier sampling



$$\hat{h} = \min_{h \in \mathcal{H}_{k, \varepsilon, M}} \|\mathbf{z} - \mathcal{A}\pi_h\|_2^2$$

GMM: result



$$x_1, \dots, x_n \stackrel{i.i.d.}{\sim} \pi_0$$

$$\omega_1, \dots, \omega_m \stackrel{i.i.d.}{\sim} \Lambda$$

$$R_{\pi_0}(h) = \mathbb{E}_{x \sim \pi_0} \ell(x, h)$$

z: Random Fourier sampling

$$\hat{h} = \min_{h \in \mathcal{H}_{k, \varepsilon, M}} \|\mathbf{z} - \mathcal{A}\pi_h\|_2^2$$

$$h^* = \min_{h \in \mathcal{H}_{k, \varepsilon, M}} R_{\pi_0}(h)$$

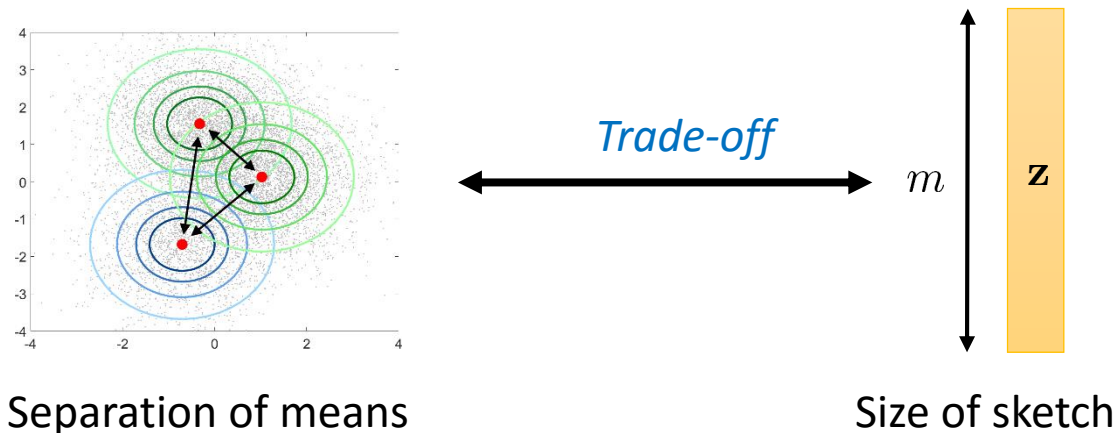
If

m large enough

Trade-off with ε_0
minimal separation

w.h.p. $R_{\pi_0}(\hat{h}) - R_{\pi_0}(h^*) \lesssim \sqrt{D_{KL}(\pi_0 \| \mathcal{H}_{k, \varepsilon, M})} + \mathcal{O}\left(\sqrt{1/n}\right)$

GMM trade-off



Separation of means	Number of measurements
$\mathcal{O}(\sqrt{d \log k})$	$m \geq \mathcal{O}(k^2 d^2 \cdot \text{polylog}(k, d))$
$\mathcal{O}(\sqrt{d + \log k})$	$m \geq \mathcal{O}(k^3 d^2 \cdot \text{polylog}(k, d))$
$\mathcal{O}(\sqrt{\log k})$	$m \geq \mathcal{O}(k^2 d^2 e^d \cdot \text{polylog}(k, d))$

More high frequencies



Sketch Size

Non-convex optimization.
Greedy heuristic: CL-OMP
[Keriven 2016]

In theory, at least

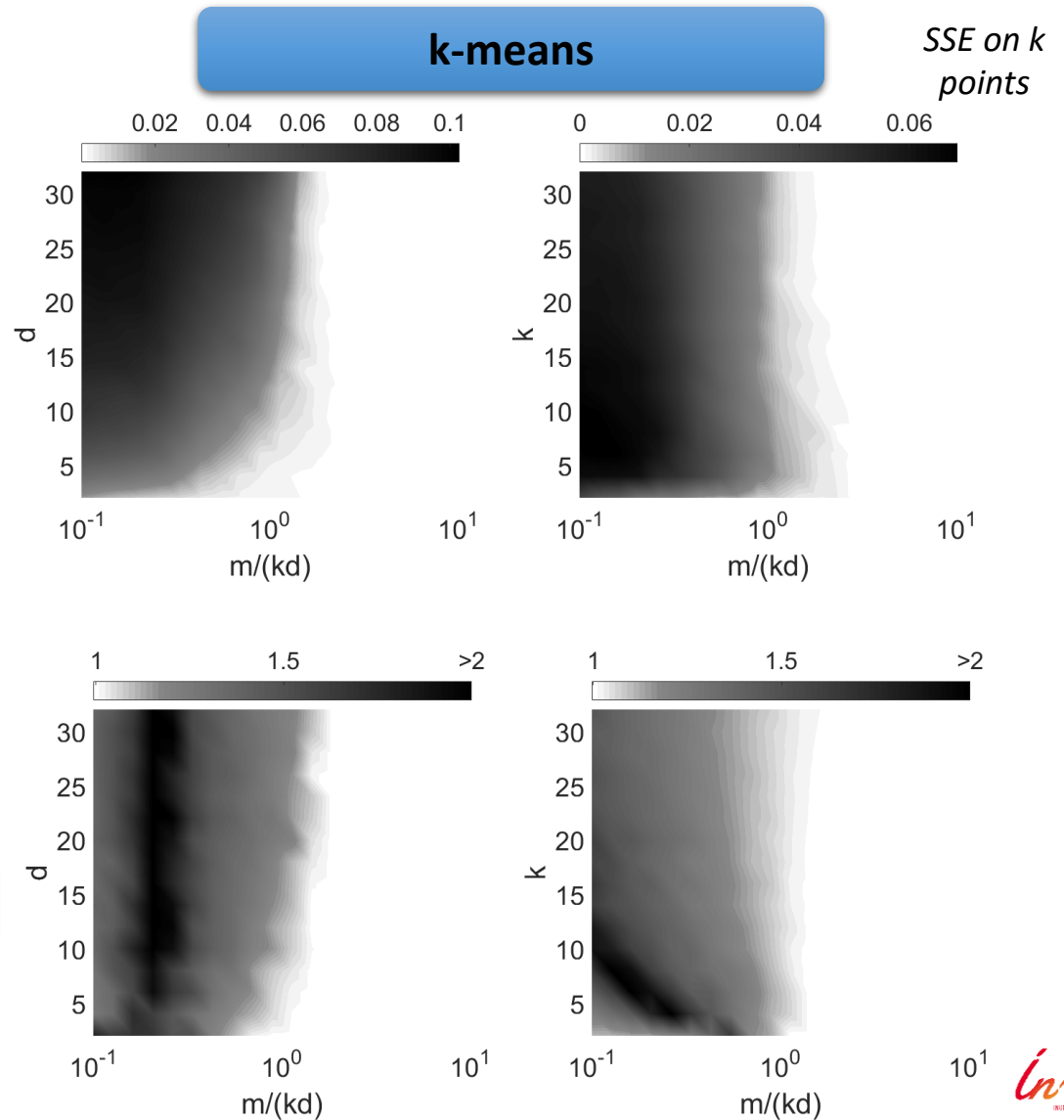
$$m \geq \mathcal{O}(k^2 d^2)$$

Empirically

$$m \approx \mathcal{O}(kd)$$

GMMs, known cov.

*Relative
loglike*



Sketch of proof

Key idea 1

Sketching operator =

Kernel mean embedding [Smola 2007]

+ Random Features [Rahimi 2007]

Step 1

Relate risk to kernel metric

Sketch of proof

Key idea 1

Sketching operator =

Kernel mean embedding [Smola 2007]

+ Random Features [Rahimi 2007]

Key idea 2

Compressive Sensing analysis

[Bourrier 2014]

Step 1

Relate risk to kernel metric

Step 2

\mathcal{A} satisfies the RIP

Sketch of proof

Key idea 1

Sketching operator =

Kernel mean embedding [Smola 2007]

+ Random Features [Rahimi 2007]

Key idea 2

Compressive Sensing analysis

[Bourrier 2014]

Main difficulty

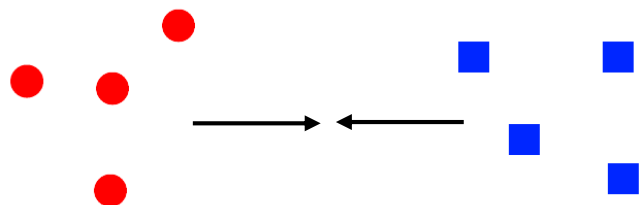
Controlling metrics between **mixtures**
that get **close to each other** in infinite-
dimensional space

Step 1

Relate risk to kernel metric

Step 2

\mathcal{A} satisfies the RIP



$\| \sum_l \alpha_l \pi_l - \sum_l \alpha'_l \pi'_l \| \rightarrow 0$: what happens ?

Sketch of proof

Key idea 1

Sketching operator =

Kernel mean embedding [Smola 2007]

+ Random Features [Rahimi 2007]

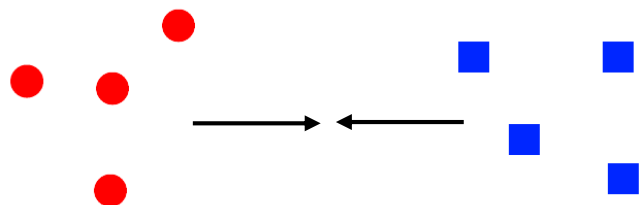
Key idea 2

Compressive Sensing analysis

[Bourrier 2014]

Main difficulty

Controlling metrics between *mixtures* that get *close to each other* in infinite-dimensional space



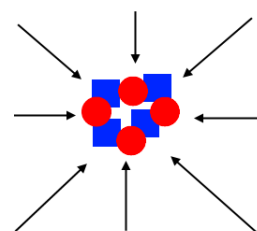
$$\left\| \sum_l \alpha_l \pi_l - \sum_l \alpha'_l \pi'_l \right\| \rightarrow 0 : \text{what happens ?}$$

Step 1

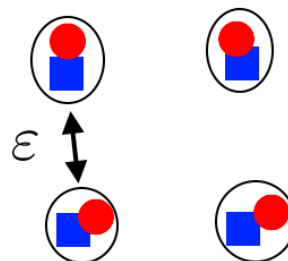
Relate risk to kernel metric

Step 2

\mathcal{A} satisfies the RIP



No hypothesis



Separation hypothesis

Outline

- ① Introduction
- ② Main results
- ③ Experimental illustration
- ④ **Conclusion**

Contributions

- Efficient **sketched mixture learning** framework, using **random generalized moments**
- Combination of many tools:
 - Kernel mean embedding
 - Random Fourier features
 - Analysis inspired by Compressive Sensing

Contributions

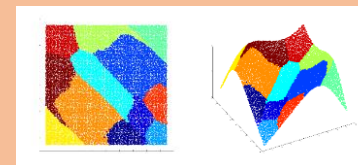
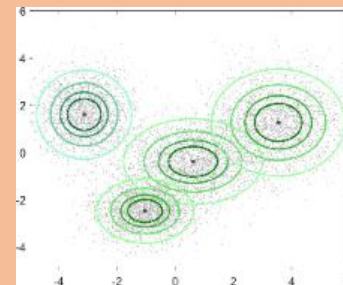
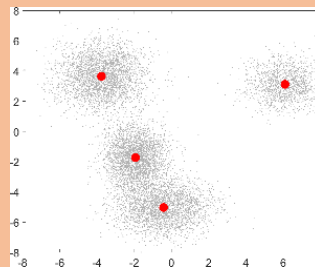
- Efficient **sketched mixture learning** framework, using **random generalized moments**
- Combination of many tools:
 - Kernel mean embedding
 - Random Fourier features
 - Analysis inspired by Compressive Sensing

Outlooks

- Bridge gap theory / practice
- Other models (done in practice), with other sketching operators
- Non-linear sketches ? (neural networks...)

SketchMLbox (sketchml.gforge.inria.fr)

- Mixture of Diracs (« K-means »)
- GMMs with known covariance
- **GMMs with unknown diagonal covariance**
- **Soon:**
 - **Mixtures of multivariate alpha-stable (only known algorithm !)**
 - Gaussian Locally Linear Mapping [Deleforge 2014]
- **Optimized for user-defined**



Thank you !

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



Appendix : CLOMPR

Algorithm 2: Compressive mixture learning *à la* OMP: CLOMP ($T = K$) and CLOMPR ($T = 2K$)

Data: Empirical sketch $\hat{\mathbf{z}}$, sketching operator \mathcal{A} , sparsity K , number of iterations $T \geq K$

Result: Support Θ , weights α

$\hat{\mathbf{r}} \leftarrow \hat{\mathbf{z}}; \Theta \leftarrow \emptyset;$

for $t \leftarrow 1$ **to** T **do**

Step 1: Find a normalized atom highly correlated with the residual with a gradient descent

$\theta \leftarrow \text{maximize}_{\theta} \left(\text{Re} \left\langle \frac{\mathcal{A}P_{\theta}}{\|\mathcal{A}P_{\theta}\|_2}, \hat{\mathbf{r}} \right\rangle_2, \text{init} = \text{rand} \right);$

end

Step 2: Expand support

$\Theta \leftarrow \Theta \cup \{\theta\};$

end

Step 3: Enforce sparsity by Hard Thresholding if needed

if $|\Theta| > K$ **then**

$\beta \leftarrow \arg \min_{\beta \geq 0} \left\| \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \beta_k \frac{\mathcal{A}P_{\theta_k}}{\|\mathcal{A}P_{\theta_k}\|_2} \right\|_2$ Select K largest entries $\beta_{i_1}, \dots, \beta_{i_K};$

 Reduce the support $\Theta \leftarrow \{\theta_{i_1}, \dots, \theta_{i_K}\};$

end

end

Step 4: Project to find weights

$\alpha \leftarrow \arg \min_{\alpha \geq 0} \left\| \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \alpha_k \mathcal{A}P_{\theta_k} \right\|_2;$

end

Step 5: Perform a gradient descent *initialized with current parameters*

$\Theta, \alpha \leftarrow \text{minimize}_{\Theta, \alpha} \left(\left\| \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \alpha_k \mathcal{A}P_{\theta_k} \right\|_2, \text{init} = (\Theta, \alpha), \text{constraint} = \{\alpha \geq 0\} \right);$

end

Update residual: $\hat{\mathbf{r}} \leftarrow \hat{\mathbf{z}} - \sum_{k=1}^{|\Theta|} \alpha_k \mathcal{A}P_{\theta_k}$

end

Normalize α such that $\sum_{k=1}^K \alpha_k = 1$