

Mixture model and sparse regularization on measure

Nicolas Jouvin (MIA PS)

Joint work with Yohann De Castro (ICJ) & Rémi Gribonval (INRIA)

Statistiques au sommet - Rochebrune 2026



Mixture models (not this one)

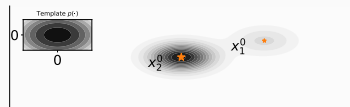
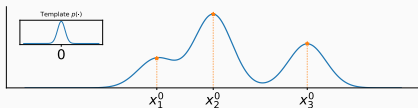


Mixture models

Context: observe sample $Z = \{z_1, \dots, z_n\} \subset \mathbb{R}^d$ from a mixture density

$$z_i \sim f^0(z) = \sum_{k=1}^s a_k^0 p(z - \mathbf{x}_k^0), \quad \mathbf{x}_k^0 \in \mathcal{X}$$

with $p \in L^1$ some known density function, e.g. Gaussian $p = \mathcal{N}(0, \Sigma)$



Goal: estimate parameters $(a_k^0, \mathbf{x}_k^0)_{k=1}^s$ of model

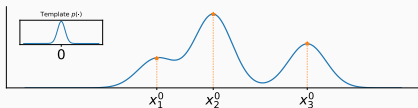
- Usually: MLE, theoretical analysis & algorithm (EM), ...

Mixture models

Context: observe sample $Z = \{z_1, \dots, z_n\} \subset \mathbb{R}^d$ from a mixture density

$$z_i \sim f^0(z) = \sum_{k=1}^s a_k^0 p(z - \mathbf{x}_k^0), \quad \mathbf{x}_k^0 \in \mathcal{X}$$

with $p \in L^1$ some known density function, e.g. Gaussian $p = \mathcal{N}(0, \Sigma)$



Goal: estimate parameters $(a_k^0, \mathbf{x}_k^0)_{k=1}^s$ of model

- ▶ Usually: MLE, theoretical analysis & algorithm (EM), ...
- ▶ **This talk:** linear inverse problem on the space of measure (De Castro et al. 2021)

Spoiler: lifting the problem to measures

Encode Target sparse measure¹ (sum of diracs)

$$\mu^0 = \sum_{k=1}^s a_k^0 \delta_{\mathbf{x}_k^0}$$

Goal recover sparse μ^0 from n -sample $Z = \{z_1, \dots, z_n\} \sim f^0 = \Phi\mu^0$

$$\text{Linear operator } \Phi : \mu \in \mathcal{M}(\mathcal{X}) \mapsto \Phi\mu := \int p(\cdot - \mathbf{x}) d\mu(\mathbf{x}) = p \star \mu$$

How ? Sparse regression on measures $\mathcal{M}(\mathcal{X})$

$$\hat{\mu} \in \arg \min_{\mu} L(Z, \mu) + R(\mu)$$

¹Bayesians call μ^0 a *mixing distribution*. When it is continuous, $\Phi\mu^0$ is a continuous mixture.

The Beurling-LASSO (BLASSO)

Reminder on the LASSO...

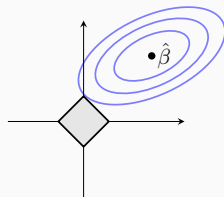
Linear combination of p covariates

$$\mathbf{y} = \mathbf{A}\boldsymbol{\beta}^0 + \boldsymbol{\Gamma}, \quad \mathbf{A} \in \mathbb{R}^{n \times p}$$

Goal: recover s -sparse $\boldsymbol{\beta}^0$ from noisy \mathbf{y}

How ? solve the LASSO problem

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \underbrace{\frac{1}{2} \|\mathbf{y} - \mathbf{A}\boldsymbol{\beta}\|^2}_{\text{data-fitting}} + \underbrace{\kappa \|\boldsymbol{\beta}\|_1}_{\text{sparsity promoting}}$$



Guarantees under conditions² on the measurement \mathbf{A}

- ▶ recovery results $\text{Supp } \hat{\boldsymbol{\beta}} = \text{Supp } \boldsymbol{\beta}^0$
- ▶ control on $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^0\|$

²irrepresentability, RIP, REV, etc.

.. and its extension to measures (BLASSO)

Linear measurements **of a measure** through

$$F : \begin{array}{l} \mathcal{M}(\mathcal{X}) \rightarrow \text{Hilbert } \mathcal{F} \\ \mu \mapsto F\mu \end{array}, \quad \text{Observe } y = F\mu_0 + \Gamma \in \mathcal{F}$$

Goal: recover s -sparse μ^0 from noisy y by solving

Beurling Lasso (De Castro and Gamboa 2012; Duval and Peyré 2015)

$$\hat{\mu} \in \arg \min_{\mu \in \mathcal{M}(\mathcal{X})} J_{\kappa}(\mu) := \frac{1}{2} \|y - F\mu\|_{\mathcal{F}}^2 + \kappa \|\mu\|_{\text{TV}} \quad (\text{BLASSO})$$

If F is continuous then there exists (at least one) $\hat{\mu}$

Total variation norm $\|\mu\|_{\text{TV}} \rightsquigarrow$ **continuous analog of l^1 -norm**

$$\text{(Discrete)} \quad \mu = \sum_{k=1}^p a_k \delta_{\mathbf{x}_k} \quad \longrightarrow \quad \|\mu\|_{\text{TV}} = \sum_{k=1}^p |a_k| = \|a\|_1$$

$$\text{(Continuous)} \quad d\mu = f d\lambda \quad \longrightarrow \quad \|\mu\|_{\text{TV}} = \|f\|_{L^1}$$

$$\hat{\mu} \in \arg \min_{\mu \in \mathcal{M}(\mathcal{X})} J_{\kappa}(\mu) := \frac{1}{2} \|y - F\mu\|_{\mathcal{F}}^2 + \kappa \|\mu\|_{TV} \quad (\text{BLASSO})$$

- ▶ Convex problem over the space of measures $\mathcal{M}(\mathcal{X})$
- ▶ How well does $\hat{\mu}$ recovers μ^0 ? Dependency on noise $\gamma = \|\Gamma\|_{\mathcal{F}}$
- ▶ Requires technical conditions on

Model kernel $\rightsquigarrow F$ induces a reproducing kernel on \mathcal{X}

$$K_{\text{mod}}(\mathbf{s}, \mathbf{t}) := \langle F\delta_{\mathbf{s}}, F\delta_{\mathbf{t}} \rangle_{\mathcal{F}}$$

with associated RKHS \mathcal{H}_{mod} which contains continuous functions

Near and far regions

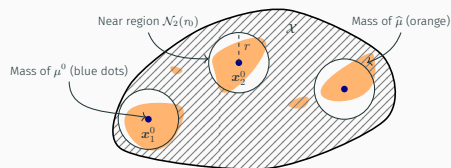
Question: how well does BLASSO estimator $\hat{\mu}$ localizes around μ^0 ?

- ▶ stated in term of near/far region around $\text{Supp } \mu^0$

Radius $r > 0$, distance $\mathfrak{d}(\cdot, \cdot)$

Near: $\mathcal{N}_k(r) := \{\mathbf{x}, \mathfrak{d}(\mathbf{x}, \mathbf{x}_k^0) \leq r\}$

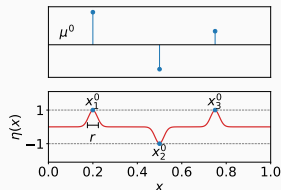
Far: $\mathcal{F}(r) := \mathcal{X} \setminus \cup_k \mathcal{N}_k(r)$



Adapted from Giard, De Castro, and Marteau (2025)

- ▶ requires on the construction of a **non-degenerate dual certificate**

1. $\eta^0 \in \mathcal{H}_{\text{mod}}$
2. $\eta^0(x_k^0) = \pm 1 \rightsquigarrow \eta^0$ in $\partial \|\mu^0\|_{\text{TV}}$
3. + additional controls on near/far



Recovery guarantees for BLASSO

Note dependency on noise & \mathcal{H}_{mod}

BLASSO recovery (informal)

Suppose $\eta^0 \in \mathcal{H}_{\text{mod}}$ is a $(r_0, \epsilon_0, \epsilon_2)$ -non-degenerate certificate. For a noise $\|\Gamma\|_{\mathcal{F}} \leq \gamma$ and regularization $\kappa \propto \frac{\gamma}{\|\eta^0\|_{\mathcal{H}_{\text{mod}}}}$ we have

1. *Small mass on far region:*

$$|\hat{\mu}|(\mathcal{F}(r_0)) \lesssim_d \gamma \|\eta^0\|_{\mathcal{H}_{\text{mod}}},$$

2. *Mass of near regions $\sim a_k^0$:*

$$|\hat{\mu}(\mathcal{N}_k(r_0)) - a_k^0| \lesssim_d \gamma \|\eta^0\|_{\mathcal{H}_{\text{mod}}},$$

3. *Detection level:* For all borelian $A \subset \mathcal{X}$ such that $|\hat{\mu}|(A) \gtrsim_d \gamma \|\eta^0\|_{\mathcal{H}_{\text{mod}}}$,

$$\exists \mathbf{x}_k^0, \quad \mathfrak{d}(A, \mathbf{x}_k^0) := \min_{t \in A} \mathfrak{d}(t, \mathbf{x}_k^0) \lesssim_d r_0,$$

Proving dual certificate is hard

State of the art analysis for support recovery (Poon, Keriven, and Peyré 2023)

Local positive curvature (LPC) assumption

If model kernel K_{mod} satisfies

[...Technical condition on K_{mod} ...]

Then there exists a $(\bar{\epsilon}_0, \bar{\epsilon}_2, r_0)$ -non-degenerate certificate.

The distance \mathfrak{d} is the *Fisher-Rao* metric associated to K_{mod} .

Proving dual certificate is hard

State of the art analysis for support recovery (Poon, Keriven, and Peyré 2023)

Local positive curvature (LPC) assumption

If model kernel K_{mod} satisfies

$$\forall 0 \leq i, j \leq 2, i + j \leq 3, B_{ij} := \sup_{\mathbf{s}, \mathbf{t} \in \mathcal{X}} \left\| K^{(i,j)}(\mathbf{s}, \mathbf{t}) \right\|_{\mathbf{s}, \mathbf{t}} < +\infty,$$

$\exists r_0 \in (0, 1/\sqrt{B_{02}})$ such that

$$\bar{\varepsilon}_0 := \sup_{\varepsilon \geq 0} \left\{ \varepsilon : K(\mathbf{s}, \mathbf{t}) \leq 1 - \varepsilon, \forall \mathbf{s}, \mathbf{t} \in \mathcal{X} \text{ s.t. } \mathfrak{d}_{\mathbf{g}}(\mathbf{s}, \mathbf{t}) \geq r_0 \right\} < +\infty$$

$$\bar{\varepsilon}_2 := \sup_{\varepsilon \geq 0} \left\{ \varepsilon : -K^{(0,2)}(\mathbf{s}, \mathbf{t})[\mathbf{v}, \mathbf{v}] \geq \varepsilon \|\mathbf{v}\|_{\mathbf{t}}^2, \forall \mathbf{v} \in \mathbb{T}_{\mathbf{t}}, \forall \mathbf{s}, \mathbf{t} \in \mathcal{X} \text{ s.t. } \mathfrak{d}_{\mathbf{g}}(\mathbf{s}, \mathbf{t}) < r_0 \right\} < +\infty$$

$$\Delta_0 < +\infty, \text{ with: } \Delta_0 = \inf \left\{ \Delta : \sum_{l=2}^s \|K^{(i,j)}(\mathbf{x}_1, \mathbf{x}_l)\|_{\mathbf{x}_1, \mathbf{x}_l} \leq \min\left(\frac{\bar{\varepsilon}_0}{B_0}, \frac{2\bar{\varepsilon}_2}{B_2}\right), i, j = 0, \dots, 2, \min_{k,l=1,\dots,s} \mathfrak{d}_{\mathbf{g}}(\mathbf{x}_k^0, \mathbf{x}_l^0) \geq \Delta \right\}$$

Then there exists a $(\bar{\varepsilon}_0, \bar{\varepsilon}_2, r_0)$ -non-degenerate certificate for Δ_0 -separated target μ^0 .

With $\mathfrak{d}_{\mathbf{g}}$ the Fisher-Rao metric associated to K_{mod} .

Proving dual certificate is hard

State of the art analysis for support recovery (Poon, Keriven, and Peyré 2023)

Local positive curvature (LPC) assumption

If model kernel K_{mod} satisfies

[...Technical condition on K_{mod} ...]

Then there exists a $(\bar{\epsilon}_0, \bar{\epsilon}_2, r_0)$ -non-degenerate certificate.

The distance \mathfrak{d} is the *Fisher-Rao* metric associated to K_{mod} .

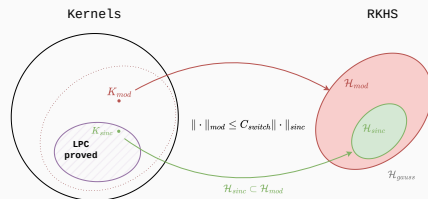
Our contribution (De Castro, Gribonval, and Jouvin 2025)

1. New measurement operator F ? You need to prove LPC for your new K_{mod} ...
 \rightsquigarrow ... or leverage on known LPC-kernel: *pivot kernels*
2. Radius r of near/far is **fixed** but in statistical applications $\gamma \sim n^{-1/2}$
 \rightsquigarrow **make the radius r_γ adaptive to noise**

The "kernel switch" principle in a nutshell

- ▶ Theory says η^0 must $\in \mathcal{H}_{\text{mod}}$
- ▶ Switch LPC kernel K_{pivot} with RKHS $\mathcal{H}_{\text{pivot}}$
- ▶ **We need**
 1. Inclusion: $\mathcal{H}_{\text{pivot}} \subset \mathcal{H}_{\text{mod}}$
 2. Control $\|\eta\|_{\mathcal{H}_{\text{mod}}} \leq C_{\text{switch}} \|\eta\|_{\mathcal{H}_{\text{pivot}}}$

The price you "pay" is C_{switch}



Explicit C_{switch} for translation-invariant $K \rightsquigarrow$ spectral representation $\nu(\omega)$

$$K(\mathbf{s}, \mathbf{t}) = \rho(\mathbf{s} - \mathbf{t}) \stackrel{\text{(Bochner)}}{=} \int e^{+i\omega^\top(\mathbf{s}-\mathbf{t})} \nu(\omega) d\omega \quad \longrightarrow \quad C_{\text{switch}}(K_1, K_2) = \inf_{\omega} \sqrt{\frac{\nu_2}{\nu_1}}(\omega)$$

But, you promised to talk about mixtures ?

Lifting on the space of measures (cont'd)

$\mathbf{a}^0, \mathbf{x}^0$ Parameter space \mathcal{X}

Measurement process



Lifting on the space of measures (cont'd)

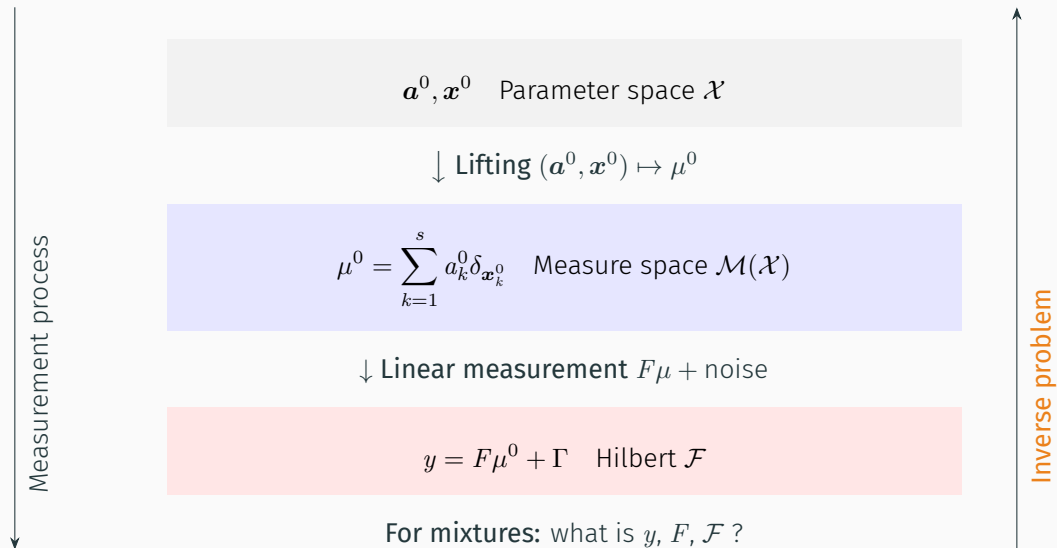
$\mathbf{a}^0, \mathbf{x}^0$ Parameter space \mathcal{X}

↓ Lifting $(\mathbf{a}^0, \mathbf{x}^0) \mapsto \mu^0$

$$\mu^0 = \sum_{k=1}^s a_k^0 \delta_{\mathbf{x}_k^0} \quad \text{Measure space } \mathcal{M}(\mathcal{X})$$

Measurement process

Lifting on the space of measures (cont'd)



Supermix: mixtures as sparse measure recovery (De Castro et al. 2021)

Observe $z_{1:n} \sim f^0$ from the mixture density. Target sparse measure

$$\mu^0 = \sum_{k=1}^s a_k^0 \delta_{x_k^0} \xrightarrow{\text{convolution}} f^0 = \Phi \mu^0 = p \star \mu^0$$

Data-fitting term ? \rightsquigarrow we only have an n -sample $z_{1:n} \sim f^0$. How to compare

$$\hat{f}_n = \frac{1}{n} \sum_{j=1}^n \delta_{z_j} \quad ? \quad \Phi \mu \quad \longrightarrow \quad \text{not in the same space}$$

Supermix: mixtures as sparse measure recovery (De Castro et al. 2021)

Observe $z_{1:n} \sim f^0$ from the mixture density. Target sparse measure

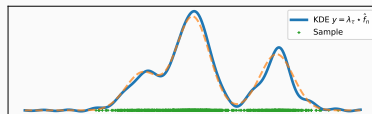
$$\mu^0 = \sum_{k=1}^s a_k^0 \delta_{x_k^0} \xrightarrow{\text{convolution}} f^0 = \Phi \mu^0 = p \star \mu^0$$

Data-fitting term ? \rightsquigarrow we only have an n -sample $z_{1:n} \sim f^0$. How to compare

$$\hat{f}_n = \frac{1}{n} \sum_{j=1}^n \delta_{z_j} \stackrel{?}{\approx} \Phi \mu \longrightarrow \text{not in the same space}$$

Embed into the same RKHS Smoothing kernel $\lambda_\tau(\cdot) = \tau^{-d} \lambda(\cdot/\tau)$

- ▶ RKHS embedding $L_\tau : \nu \mapsto \lambda_\tau \star \nu \in \mathcal{F}_\tau$
- ▶ Associated Hilbert \mathcal{F}_τ depends on λ_τ
- ▶ Observation is the KDE
 $y = L_\tau \hat{f}_n = \lambda_\tau \star \hat{f}_n$



KDE using sinus cardinal $\lambda(x) = \sin(x)/x$

Solve the following BLASSO problem

$$\hat{\mu} \in \arg \min_{\mu} \frac{1}{2} \left\| \underbrace{L_{\tau} \hat{f}_n}_y - \underbrace{(L_{\tau} \circ \Phi) \mu}_{F_{\tau}} \right\|_{\mathcal{F}_{\tau}}^2 + \kappa \|\mu\|_{\text{TV}} \quad (\text{Supermix})$$

Model kernel is explicit & translation-invariant $K_{\text{mod}} = \rho_{\text{mod}} = \lambda_{\tau} \star p \star \check{p}$

However

1. *Theory*: LPC needs to be checked for each template $p \rightsquigarrow$ **switch!**

$$C_{\text{switch}} = C_{\text{switch}}(\tau, p) := \sqrt{\frac{\tau^d}{\inf_{\mathbb{B}_{\tau}} \nu_{\lambda_{\tau}} |\nu_p(\omega)|^2}} < +\infty, \quad (\mathbf{H}_p)$$

Solve the following BLASSO problem

$$\hat{\mu} \in \arg \min_{\mu} \frac{1}{2} \left\| \underbrace{L_{\tau} \hat{f}_n}_y - \underbrace{(L_{\tau} \circ \Phi) \mu}_{F_{\tau}} \right\|_{\mathcal{F}_{\tau}}^2 + \kappa \|\mu\|_{\text{TV}} \quad (\text{Supermix})$$

Model kernel is explicit & translation-invariant $K_{\text{mod}} = \rho_{\text{mod}} = \lambda_{\tau} \star p \star \check{p}$

However

1. *Theory*: LPC needs to be checked for each template $p \rightsquigarrow$ **switch!**

$$C_{\text{switch}} = C_{\text{switch}}(\tau, p) := \sqrt{\frac{\tau^d}{\inf_{\mathbb{B}_{\tau}} \nu_{\lambda_{\tau}} |\nu_p(\omega)|^2}} < +\infty, \quad (\mathbf{H}_p)$$

2. *Practice*: Every algorithm computing $\hat{\mu}$ requires evaluation of K_{mod} (and ∇K_{mod})

$$\text{Costly } d\text{-dimensional integrals } K_{\text{mod}}(\mathbf{s} - \mathbf{t}) = \int_{\mathbb{R}^d} (\dots)$$

Sketching the kernel: random Fourier features (Rahimi and Recht 2007)

K_{mod} admits a weighted *random Fourier features* representation

$$K_{\text{mod}}(\mathbf{s}, \mathbf{t}) = \mathbb{E}_{\omega \sim \Lambda} \left[\varphi_{\omega}(\mathbf{s}) \overline{\varphi_{\omega}(\mathbf{t})} \right], \quad \varphi_{\omega}(\mathbf{t}) := W(\omega) e^{-i\omega^T \mathbf{t}}$$

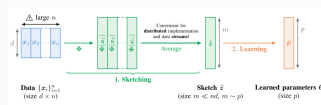
$$W(\omega) = \sqrt{\frac{\nu_{\lambda_{\tau}}}{\Lambda}} \mathfrak{F}[p](\omega)$$

- ▶ **Sketching** \rightsquigarrow Monte-Carlo approximation draw $\omega_{1:m} \sim \Lambda$
- ▶ *Sketched model kernel*

$$K_{\text{sketch,mod}}(\mathbf{s}, \mathbf{t}) := \frac{1}{m} \sum_{i=1}^m \varphi_{\omega_i}(\mathbf{s}) \overline{\varphi_{\omega_i}(\mathbf{t})} \xrightarrow[m \rightarrow +\infty]{\text{a.s.}} K_{\text{mod}}$$

- ▶ Sketch vector = compressed dataset (Gribonval et al. 2020)

$$\mathbf{y}_{\text{sketch}} = \frac{1}{\sqrt{m}} \left(\sqrt{\frac{\nu_{\lambda_{\tau}}}{\Lambda}} \mathfrak{F}[\hat{f}_n](\omega_i) \right)_{i=1}^m = \frac{1}{n} \sum_{j=1}^n \underbrace{\xi(z_j)}_{\in \mathbb{C}^m}$$



Sketching as dataset compression from Gribonval et al. (2020)

$$\hat{\mu}_{\text{sketch}} \in \arg \min_{\mu} \frac{1}{2} \|\mathbf{y}_{\text{sketch}} - F_{\text{sketch}}\mu\|_{\mathbb{C}^m}^2 + \kappa \|\mu\|_{\text{TV}} \quad (\text{S2mix})$$

Random measurements m random draws $\omega_{1:m}$ with sketched operator:

$$F_{\text{sketch}}\mu := \frac{1}{\sqrt{m}} \left(\int_{\mathcal{X}} \varphi_{\omega_i}(\mathbf{t}) d\mu(\mathbf{t}) \right)_{i=1}^m \in \mathbb{C}^m =: \mathcal{F}_{\text{sketch}}$$

$$\hat{\mu}_{\text{sketch}} \in \arg \min_{\mu} \frac{1}{2} \|\mathbf{y}_{\text{sketch}} - F_{\text{sketch}}\mu\|_{\mathbb{C}^m}^2 + \kappa \|\mu\|_{\text{TV}} \quad (\text{S2mix})$$

Random measurements m random draws $\omega_{1:m}$ with sketched operator:

$$F_{\text{sketch}}\mu := \frac{1}{\sqrt{m}} \left(\int_{\mathcal{X}} \varphi_{\omega_i}(\mathbf{t}) d\mu(\mathbf{t}) \right)_{i=1}^m \in \mathbb{C}^m =: \mathcal{F}_{\text{sketch}}$$

Question How many measurement m to ensure recovery with high \mathbb{P}_{Λ} ?

Informal: De Castro, Gribonval, and Jouvin (2025, Proposition 3.2)

If the sketch size

$$m \gtrsim \text{cte}(d, \tau) \cdot s \cdot \log \left(\frac{s}{\alpha} \right)$$

Then, with proba $1 - \alpha$, the recovery guarantee hold for S2Mix with kernel switch

Effective near regions: beyond fixed radius

What about noise ? We can show $\gamma_n = \|y - F\mu^0\|_{\mathcal{F}} \leq \frac{1}{\sqrt{n}}$

Effective near region guarantees hold for radius $r_n = n^{-1/4}\delta_n$, with any $\delta_n \rightarrow +\infty$

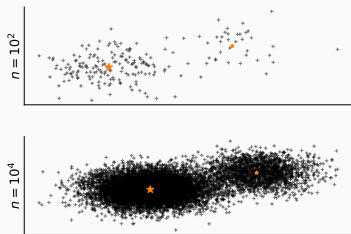
$$n \gg 1, \quad \hat{\mu}(\mathcal{F}(r_n)) \lesssim \frac{1}{\delta_n^2} C_{\text{switch}} \|\eta^0\|_{\mathcal{H}_{\text{pivot}}}$$

Effective near regions: beyond fixed radius

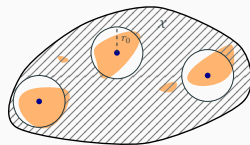
What about noise ? We can show $\gamma_n = \|y - F\mu^0\|_{\mathcal{F}} \leq \frac{1}{\sqrt{n}}$

Effective near region guarantees hold for radius $r_n = n^{-1/4}\delta_n$, with any $\delta_n \rightarrow +\infty$

$$n \gg 1, \quad \hat{\mu}(\mathcal{F}(r_n)) \lesssim \frac{1}{\delta_n^2} C_{\text{switch}} \|\eta^0\|_{\mathcal{H}_{\text{pivot}}}$$



Same radius r_0 →

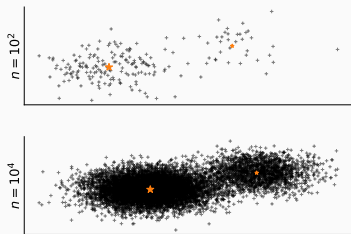


Effective near regions: beyond fixed radius

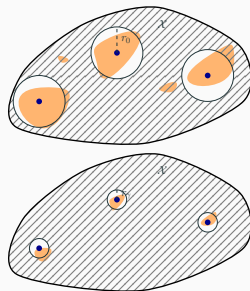
What about noise ? We can show $\gamma_n = \|y - F\mu^0\|_{\mathcal{F}} \leq \frac{1}{\sqrt{n}}$

Effective near region guarantees hold for radius $r_n = n^{-1/4}\delta_n$, with any $\delta_n \rightarrow +\infty$

$$n \gg 1, \quad \hat{\mu}(\mathcal{F}(r_n)) \lesssim \frac{1}{\delta_n^2} C_{\text{switch}} \|\eta^0\|_{\mathcal{H}_{\text{pivot}}}$$



Noise dependent r_{γ} \rightarrow



Conclusion

Contributions

- ▶ **Kernel switch:** general principle to build certificate without proving LPC
- ▶ **Effective near regions** noise adaptive radius r_γ





Perspective

- ▶ Beyond TI mixtures \rightsquigarrow GMM $\mathbf{x}_k^0 = (\mu_k, \Sigma_k)$ (Giard, De Castro, and Marteau 2025)
- ▶ *Computational side:* practical algorithms to compute $\hat{\mu}$ (Chizat 2019; De Castro, Gadat, and Marteau 2023)

Gradient descent for $\min_{\text{discrete } \mu_p} J_\kappa(\mu_p) \xrightarrow[p \rightarrow +\infty]{\text{approximate Wasserstein flow}} \frac{\partial \mu_t}{\partial t} = -\text{div}(\mu_t \nabla^W J_\kappa(\mu_t))$

Thank you for your attention !

References

-  Chizat, Lenaic (2019). **“Sparse optimization on measures with over-parameterized gradient descent”**. In: *arXiv preprint arXiv:1907.10300*.
-  De Castro, Yohann, Sébastien Gadat, and Clément Marteau (2023). **“FastPart: Over-Parameterized Stochastic Gradient Descent for Sparse optimisation on Measures”**. In: *arXiv preprint arXiv:2312.05993*.
-  De Castro, Yohann and Fabrice Gamboa (2012). **“Exact reconstruction using Beurling minimal extrapolation”**. In: *Journal of Mathematical Analysis and applications* 395.1, pp. 336–354.
-  De Castro, Yohann, Rémi Gribonval, and Nicolas Jouvin (2025). **“Effective regions and kernels in continuous sparse regularisation, with application to sketched mixtures”**. In: *arXiv preprint arXiv:2507.08444*.

Bibliography ii

-  De Castro, Yohann et al. (2021). **“SuperMix: Sparse regularization for mixtures”**. In: *The Annals of Statistics* 49.3, pp. 1779 –1809.
-  Duval, V. and G. Peyré (2015). **“Exact support recovery for sparse spikes deconvolution”**. In: *Foundations of Computational Mathematics*, pp. 1–41.
-  Giard, Romane, Yohann De Castro, and Clément Marteau (2025). **“Gaussian Mixture Model with unknown diagonal covariances via continuous sparse regularization”**. In: *arXiv preprint arXiv:2509.12889*.
-  Gribonval, Rémi et al. (2020). **“Sketching datasets for large-scale learning (long version)”**. In: *arXiv preprint arXiv:2008.01839*.
-  Poon, Clarice, Nicolas Keriven, and Gabriel Peyré (2023). **“The geometry of off-the-grid compressed sensing”**. In: *Foundations of Computational Mathematics* 23.1, pp. 241–327.
-  Rahimi, Ali and Benjamin Recht (2007). **“Random features for large-scale kernel machines”**. In: *Advances in neural information processing systems* 20.

Why is it called a *dual* certificate ? Duality

Convex program The BLASSO problem admit a (pre)-dual

$$\inf_{\mu \in \mathcal{M}(\mathcal{X})} \frac{1}{2} \|y - F\mu\|_{\mathcal{F}}^2 + \kappa \|\mu\|_{\text{TV}} \xrightarrow{\text{Dual}} \inf_{c \in \mathcal{F} : \|F^*c\|_{\infty} \leq 1} \langle c, y \rangle_{\mathcal{F}} - \frac{\kappa}{2} \|c\|_{\mathcal{F}}^2 \quad (\mathcal{D}_{\kappa})$$

With the adjoint $F^* : \mathcal{F} \rightarrow \mathcal{H}_{\text{mod}}$

Optimality conditions yield for a BLASSO minimizer $\hat{\mu}$

$$0 \in \partial J_{\kappa}(\hat{\mu}) \iff 0 \in \partial \|\hat{\mu}\|_{\text{TV}} + \frac{1}{\kappa} F^*(y - F\hat{\mu}) \iff \eta := -\frac{1}{\kappa} F^*(y - F\hat{\mu}) \in \partial \|\hat{\mu}\|_{\text{TV}}$$

We have

- ▶ $\text{Supp } \hat{\mu} \subset \{\mathbf{x} \in \mathcal{X}, |\eta(\mathbf{x})| = 1\} \longrightarrow \eta$ certifies the support of $\hat{\mu}$
- ▶ Natural to look for certificates $\eta^0 \in \text{Im}(F^*)$ s.t. $\eta^0 \in \partial \|\mu^0\|_{\text{TV}}$

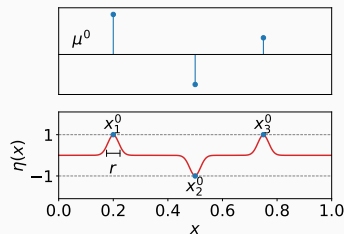
Non-degeneracy condition for the dual certificate

Find η^0 subgradient of $\|\mu^0\|_{\text{TV}}$ with additional controls on near/far

Fix radius $r > 0$, distance $\mathfrak{d}(\cdot, \cdot)$ & control $\epsilon_0, \epsilon_2 > 0$.

η^0 is a $(r, \epsilon_0, \epsilon_2)$ -**non-degenerate** certificate iff

- ▶ $\eta^0 \in \text{Im}(F^*) = \mathcal{H}_{\text{mod}}$
- ▶ $\eta^0(\mathbf{x}_k^0) = \text{sign}(a_k^0)$
- ▶ $\eta^0(\mathbf{x}) \leq 1 - \epsilon_0, \forall \mathbf{x} \in \mathcal{F}(r)$
- ▶ $\eta^0(\mathbf{x}) \leq 1 - \epsilon_2 \mathfrak{d}(\mathbf{x}, \mathbf{x}_k^0)^2, \forall \mathbf{x} \in \mathcal{N}_k(r)$



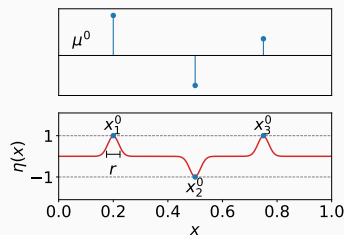
Non-degeneracy condition for the dual certificate

Find η^0 subgradient of $\|\mu^0\|_{\text{TV}}$ with additional controls on near/far

Fix radius $r > 0$, distance $\mathfrak{d}(\cdot, \cdot)$ & control $\epsilon_0, \epsilon_2 > 0$.

η^0 is a $(r, \epsilon_0, \epsilon_2)$ -**non-degenerate** certificate iff

- ▶ $\eta^0 \in \text{Im}(F^*) = \mathcal{H}_{\text{mod}}$
- ▶ $\eta^0(\mathbf{x}_k^0) = \text{sign}(a_k^0)$
- ▶ $\eta^0(\mathbf{x}) \leq 1 - \epsilon_0, \forall \mathbf{x} \in \mathcal{F}(r)$
- ▶ $\eta^0(\mathbf{x}) \leq 1 - \epsilon_2 \mathfrak{d}(\mathbf{x}, \mathbf{x}_k^0)^2, \forall \mathbf{x} \in \mathcal{N}_k(r)$



Certificates are the main workhorse: \exists NDC \implies recovery guarantees

Candidate certificate: For a kernel K satisfying LPC build

$$\eta^0 = \sum_{k=1}^s \alpha_k K(\mathbf{x}_k^0, \cdot) + \sum_{k=1}^s \langle \beta_k, \nabla_1 K(\mathbf{x}_k^0, \cdot) \rangle \in \mathcal{H}_K$$

\rightsquigarrow linear system in $(\alpha_k, \beta_k)_{k=1}^s$ for the constraints $\eta^0(\mathbf{x}_k^0) = 1$ and $\nabla \eta^0(\mathbf{x}_k^0) = 0$

- ▶ *Invertibility?* The system is invertible if spikes are “sufficiently separated”
- ▶ *Non-degeneracy* $\iff K$ satisfies LPC
- ▶ Under LPC we also have $\|\eta^0\|_{\mathcal{H}_K} \leq \sqrt{s}$

Different options

1. Prove K_{mod} satisfy LPC + $K = K_{\text{mod}}$ for certificate \rightsquigarrow cumbersome
2. *Kernel switch:* use a pivot kernel $K = K_{\text{pivot}} \longrightarrow \eta^0 \in \mathcal{H}_{K_{\text{pivot}}} \subset \mathcal{H}_{\text{mod}}$

From Bregman to near/far region

We can always bound the Bregman divergence

$$\begin{aligned} D_{\eta^0}(\widehat{\mu} \parallel \mu^0) &:= \|\widehat{\mu}\|_{\text{TV}} - \|\mu^0\|_{\text{TV}} - \int \eta^0 d(\widehat{\mu} - \mu^0) \\ &\leq (\dots) \leq \frac{(\gamma + \kappa \|\eta^0\|_{\mathcal{H}_{\text{mod}}})}{2\kappa} \end{aligned}$$

Choosing $\kappa = \gamma / \|\eta^0\|_{\mathcal{H}_{\text{mod}}} +$ using the control of η^0 we get for any radius $r \leq r_0$,

$$\begin{aligned} \mathcal{D}_{\eta^0}(\mu \parallel \mu^0) &= \|\mu\|_{\text{TV}} - \langle \eta^0, \mu \rangle_{\mathcal{C}(\mathcal{X}) \times \mathcal{M}(\mathcal{X})} \\ &\geq \|\mu\|_{\text{TV}} - \sum_{l=1}^s \int_{\mathcal{N}_l(r)} |\eta^0| d|\mu| - \int_{\mathcal{F}(r)} |\eta^0| d|\mu| \end{aligned}$$

...

$$\mathcal{D}_{\eta^0}(\mu \parallel \mu^0) \geq \bar{\varepsilon}_2 r^2 |\mu|(\mathcal{F}(r)) + \bar{\varepsilon}_2 \sum_{l=1}^s \int_{\mathcal{N}_l(r)} \mathfrak{d}_{\mathbf{g}}(x, t_i^0)^2 d|\mu|(x). \quad (1)$$

Sketching

Sketch vector = compressed dataset

$$\mathbf{y}_{\text{sketch}} = \frac{1}{\sqrt{m}} \left(\sqrt{\frac{\nu \lambda_\tau}{\Lambda}} \mathfrak{F}[\hat{f}_n](\omega_i) \right)_{i=1}^m$$

Analogies with compressed sensing

$$\mathbf{y}_{\text{sketch}} = \frac{1}{n} \sum_{j=1}^n \underbrace{\xi(\mathbf{z}_j)}_{\in \mathbb{C}^m} \xrightarrow[n \rightarrow +\infty]{\text{a.s.}} \int \xi(\mathbf{z}) f^0(\mathbf{z}) d\mathbf{z} = F_{\text{sketch}} \mu^0$$

with $\xi(\mathbf{z}) \in \mathbb{C}^m$ the (weighted) random Fourier measurements of the dataset

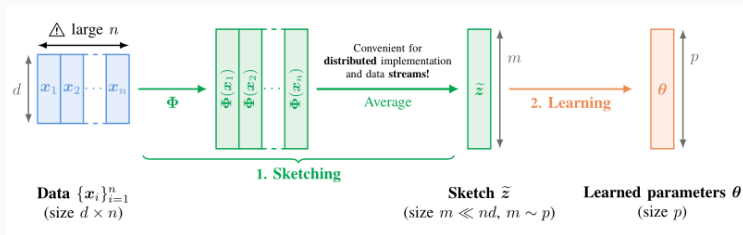


Illustration of sketching as dataset compression from Gribonval et al. (2020)

Model RKHS & functional analysis

$$\begin{aligned}K_{\text{mod}}(\mathbf{s}, \mathbf{t}) &:= \langle F\delta_{\mathbf{s}}, F\delta_{\mathbf{t}} \rangle_{\mathcal{F}}, \\ &= \langle K_{\text{mod}}(\mathbf{s}, \cdot), K_{\text{mod}}(\mathbf{t}, \cdot) \rangle_{\mathcal{H}_{\text{mod}}}\end{aligned}$$

Adjoint $F^* : \mathcal{F} \rightarrow \mathcal{H}_{\text{mod}}$ and $(F^*c)(\mathbf{t}) = \eta_c(\mathbf{t}) := \langle c, F\delta_{\mathbf{t}} \rangle_{\mathcal{F}}$

- ▶ The unique RKHS of K_{mod} is given by

$$\mathcal{H}_{\text{mod}} = \{ \eta : \mathcal{X} \rightarrow \mathbb{R} \mid \exists c \in \mathcal{F}, \eta = \eta_c \} = \{ \eta : \mathcal{X} \rightarrow \mathbb{R} \mid \exists ! c \in \overline{\text{Im}(F)}, \eta = \eta_c \},$$

and for all c orthogonal to $\overline{\text{Im}(F)}$ in \mathcal{F} , $\eta_c = 0$.

- ▶ The isometry is given by the mapping $c \in \overline{\text{Im}(F)} \mapsto \eta_c \in \mathcal{H}_{\text{mod}}$.
- ▶ The norms satisfy

$$\forall c \in \overline{\text{Im}(F)}, \quad \|\eta_c\|_{\mathcal{H}_{\text{mod}}} = \|c\|_{\mathcal{F}}.$$

Mixture models: the kernel switch constant

the case of *supersmooth* densities

$$\exists p \in [1, +\infty], \alpha, \beta > 0, \quad \mathfrak{F}[p](\boldsymbol{\omega}) \propto e^{-\alpha \|\boldsymbol{\omega}\|_p^\beta}$$

leads to a scaling in

$$C_{\text{switch}} = \mathcal{O}_d \left(\tau^{d/2} e^{\alpha \left(\frac{d^{1/p}}{\tau} \right)^\beta} \right),$$

\mathcal{O}_d may depend exponentially in the dimension d but not on τ . This case encompasses

- ▶ Gaussian $\rightsquigarrow C_{\text{switch}} = \mathcal{O}_d(\tau^{d/2} e^{d/2\tau^2})$
- ▶ multivariate Cauchy
- ▶ product of univariate Cauchy
- ▶ centered stable distributions with known scale parameter and zero skewness.