

# *Bayesian sparse optimization for large scale dictionary learning*

Daniela Calvetti

Case Western Reserve University

Department of Mathematics, Applied Mathematics and Statistics

Based on joint work with Erkki Somersalo, Nathan Waniorek and Alberto Bocchinfuso.

Supported in part by NSF DMS-1951446

AWM Workshop: Women in Inverse Problems, June 1, 2023

# Dictionary matching formulation

Consider the system:

$$Dx = b, \quad D \in \mathbb{R}^{m \times n}, \quad b \in \mathbb{R}^m,$$

where

- $D$  is a dictionary of  $n$  atoms.
- $b$  is the datum to represent .
- $x_j$  is the weight of the contribution of  $j$ -th atom to  $b$ .

# Sparse coding

Assume  $b$  admits a sparse dictionary representation, and we want to find  $\hat{x}$

$$\hat{x} = \arg \min_x \|Dx - b\|, \quad \|x\|_0 \ll n$$

To reduce complexity, replace  $D$  with a low rank factorization

$$D = WH + E, \quad W \in \mathbb{R}^{m \times k}, \quad H \in \mathbb{R}^{k \times n}, \quad k \ll n,$$

then solve

$$b = Wh + \epsilon', \quad h \geq 0,$$

where  $\epsilon'$  accounts for all discrepancy contributions.

# A multiphase approach

To reduce complexity and promote sparsity we suggest a four step approach.

- 1 **Clustering**: group original atoms into subdictionaries  $D^{(j)}$ ,  $1 \leq j \leq K$ .
- 2 **Reduction**: compress each subdictionary into low rank code

$$D^{(j)} = W^{(j)}H^{(j)} + E^{(j)}.$$

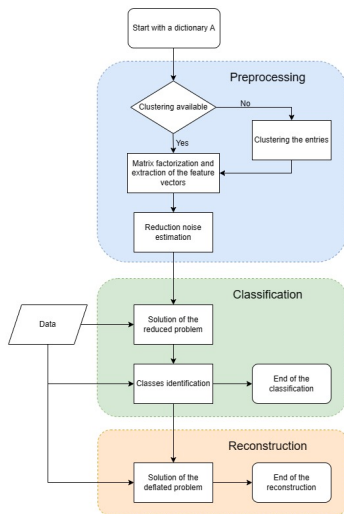
- 3 **Selection**: Identify clusters needed to explain  $b$ ,

$$b = W^{(j_1)}h^{(j_1)} + \dots + W^{(j_r)}h^{(j_r)} + \epsilon'.$$

- 4 **Deflate and Solve**: Express  $b$  in terms of identified subdictionaries,

$$b = D^{(j_1)}x^{(j_1)} + \dots + D^{(j_r)}x^{(j_r)} + \tilde{\epsilon}'.$$

# Flowchart summary



# Clustering

- Applied unsupervised clustering to dictionary atoms, to exploit latent structure of data.
- Use k-medoids.
- Regroup atoms into block-wise dictionaries.

$$D = [ D^{(1)} \mid \dots \mid D^{(k)} ]$$

# Reduction

Compute low-rank approximation of each subdictionary

$$D^{(j)} = W^{(j)}H^{(j)} + E^{(j)}, \quad W^{(j)} \in \mathbb{R}^{m \times k_j}, \quad H^{(j)} \in \mathbb{R}^{k_j \times n}, \quad k_j < n.$$

where

- $H^{(j)} \geq 0$ ;
- Feature vector = columns of  $W^{(j)}$ , retain characteristics of atoms in  $D^{(j)}$ .

# Reduction

To find  $W$ ,  $H$ , proceed as in NMF, relaxing nonnegativity for  $W$ .

- Given  $W$ , columns of  $H$  satisfy

$$h^\ell = \arg \min_h \|d^{(\ell)} - Wh\|_2$$

- Given  $H$ ,  $W$  satisfies

$$\begin{aligned} W &= \arg \min_V \|D - VH\|_F \\ &= \arg \min_V \|D^T - H^T V^T\|_F \end{aligned}$$



# NNLS- NMF algorithm

**Given:** data matrix  $D \in \mathbb{R}^{m \times n}$ ,  $X \geq 0$ , rank  $k > 0$ ,

**Initialize:**  $W^0 \in \mathbb{R}^{m \times k}$ ,  $t = 0$ .

**Iterate** until stopping criterion is met:

- ① Update H: Set  $h^{(j)} = \arg \min_{h \geq 0} \|d^{(j)} - W^{(t)}h\|_2$  for  $1 \leq j \leq n$ , and set

$$H^{(t+1)} = [ h^{(1)} \quad \dots \quad h^{(n)} ];$$

- ② Scale the rows of  $H^{(t+1)}$ ,

$$h_{(j)} \rightarrow \frac{h_{(j)}}{\|h_{(j)}\|_1}, \quad 1 \leq j \leq k.$$

- ③ Update W: Set  $w_{(j)} = \arg \min_{w \geq 0} \| (H^{(t+1)})^T w - d_{(j)} \|_2$  for  $1 \leq j \leq m$ , and set

$$W^{(t+1)} = \begin{bmatrix} w_{(1)}^T \\ \vdots \\ w_{(m)}^T \end{bmatrix};$$

- ④ Check the convergence criterion, advance the counter  $t \rightarrow t + 1$ .

# Dictionary Reduction

Approximate the dictionary matrix  $D$  by:

$$\begin{aligned}
 [ D^{(1)} \mid \dots \mid D^{(k)} ] &\approx [ W^{(1)} \mid \dots \mid W^{(k)} ] \begin{bmatrix} H^{(1)} & 0 & \dots & 0 \\ 0 & H^{(2)} & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \dots & 0 & H^{(k)} \end{bmatrix} \\
 &= [ W^{(1)}H^{(1)} \mid \dots \mid W^{(k)}H^{(k)} ] + E
 \end{aligned}$$

then

$$b = W \underbrace{Hz}_{=h} = W \begin{bmatrix} H^{(1)}z^{(1)} \\ \vdots \\ H^{(k)}z^{(k)} \end{bmatrix} + \epsilon.$$

# Dictionary Compression Error

Want to analyze

$$E = D - WH.$$

Observe that

$$\begin{aligned} d^{(\ell)} = \text{De}_\ell &= Wh^{(\ell)} + (\text{De}_\ell - Wh^{(\ell)}) \\ &= Wh^{(\ell)} + m^{(\ell)}, \quad 1 \leq \ell \leq n. \end{aligned}$$

where  $m^{(\ell)}$  is the dictionary compression error (DCE). In Bayesian framework, model  $m^{(\ell)}$  as a **random variable**. To estimate the probability density of DCE

- Regard atom  $d^{(\ell)}$  as a realization from underlying distribution  $\pi_d$
- $\pi_{\text{DCE}}$  is the **push forward** of  $\pi_d$  by mapping  $d^{(\ell)} \rightarrow m^{(\ell)}$ .
- Note that  $h^{(\ell)} = \arg \min_h \|d^{(\ell)} - Wh\|_2^2$ .
- It is not clear that discrepancies  $m^{(\ell)}$  are scaled white noise.

# Model discrepancy as noise

## Dictionary Compression Error

$$m = b - WHx$$

To estimate  $m$ ,

- 1 Draw  $s$  uniformly distributed random realizations,  $\tilde{x}^1, \dots, \tilde{x}^m$ , of  $x$ ,
- 2 Compute  $\mu_{\text{DCE}} = \frac{1}{m} \sum_{j=1}^m (b - WH\tilde{x}^j)$ .
- 3 Compute  $C_{\text{DCE}} = \frac{1}{m} \sum_{j=1}^m (b - WH\tilde{x}^j - \mu)(b - WH\tilde{x}^j - \mu)^T + \epsilon I$ .
- 4 Use the Laplace approximation  $\epsilon \sim \mathcal{N}(\mu_{\text{DCE}}, C_{\text{DCE}})$ .

Whitening

$$C_{\text{DCE}}^{-1/2}(b - \mu_{\text{DCE}}) = C_{\text{DCE}}^{-1/2}Wh + w$$

where  $w$  is white noise.

# Cluster detection

Identify subdictionaries contributing to  $b$ .

If only few subdictionaries contribute, solve

$$h = \arg \min \left\{ \left\| C_{\text{DCE}}^{-1/2} (b - \mu_{\text{DCE}} - Wh) \right\|^2 + \lambda \left\| (\|h^{(1)}\|_{(1)}, \dots, \|h^{(K)}\|_{(K)}) \right\|_0 \right\},$$

where the norms  $\| \cdot \|_{(j)}$ , are constructed by using the full dictionary information.

# Structural prior for subcluster coefficients

- The  $h^{(j)}$  is expected to have a structural similarity to the column vectors of the corresponding matrices  $H^{(j)}$  arising from the full dictionary  $D^{(j)}$ .
- $h^{(j)}$  is not expected to be close to the group mean, as most groups  $W^{(j)}$  contribute little or not at all to explain  $b$ .
- $h^{(j)}$  should be in the general direction of  $H^{(j)}$ , with amplitude controlled by a group sparsity promoting prior.

# Structural prior

Let

$$H^{(j)} = U^{(j)} \Sigma^{(j)} [V^{(j)}]^\top, \quad \Sigma^{(j)} = \text{diag}([\sigma_1^{(j)}, \dots, \sigma_{r_j}^{(j)}, 0, \dots, 0]),$$

and define

$$\begin{aligned} C^{(j)} &= \left( \frac{1}{\sigma_1^{(j)}} \right)^2 U^{(j)} \Sigma^{(j)} [\Sigma^{(j)}]^\top [U^{(j)}]^\top + \epsilon I \\ &= u_1^{(j)} [u_1^{(j)}]^\top + \sum_{\ell=2}^{r_j} \left( \frac{\sigma_\ell^{(j)}}{\sigma_1^{(j)}} \right)^2 u_\ell^{(j)} [u_\ell^{(j)}]^\top + \epsilon I, \end{aligned}$$

The *direction vector*  $\hat{h}^{(j)}$  of  $h^{(j)}$ , follows the angular central Gaussian measure with covariance  $C^{(j)}$ , defined as the pullback of the Gaussian measure  $\mathcal{N}(0, C^{(j)})$  with respect to the radial projection  $\mathbb{R}^{k_j} \rightarrow \mathbb{S}^{k_j-1}$ .

Alternatively, we may express this prior in terms of the conditional distribution, given the length of the vector, as

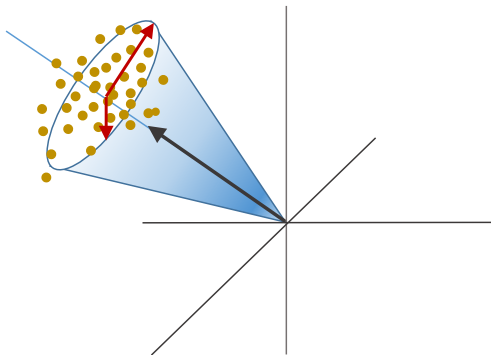
$$\pi^{(j)}(\hat{h}^{(j)} \mid \|h^{(j)}\|) \propto \exp(-\|h^{(j)}\|_{(j)}^2),$$

where the **structural norm** is defined as

$$\|h^{(j)}\|_{(j)}^2 = [h^{(j)}]^\top [C^{(j)}]^{-1} h^{(j)}. \quad (1)$$

This expresses belief that  $h^{(j)}$  could have an arbitrarily large component along  $u_1^{(j)}$ , but it should remain within the cone defined by the columns of  $H^{(j)}$ .





# Thresholding

- After computing  $h^{(j)}$  with group sparsity, cannot expect most  $\|h^{(j)}\|_{(j)}^2$  to be zero.
- We expect many  $h^{(j)}$  to be negligible.
- The only significant  $h^{(j)}$  are those above std of DCE, that is,

$$h^{(j)} \text{ is significant if and only if } \|W^{(j)} h^{(j)}\|^2 > \text{trace}(C_{\text{DCE}}^{(j)}).$$

- Let  $J = \{j_1, \dots, j_r\} \subset \{1, 2, \dots, K\}$  be the indices of relevant subdictionary.

# Deflation

Once identified the relevant classes:

- Deflate original dictionary  $D$  removing non contributing classes and solve

$$\left[ D^{(i_1)} \mid \dots \mid D^{(i_\ell)} \right] \begin{bmatrix} x^{(i_1)} \\ \vdots \\ x^{(i_\ell)} \end{bmatrix} = \hat{D}\hat{x} = b.$$

- If all the relevant classes have been identified, then

$$Dx = \hat{D}\hat{x} = b.$$

- $b$  admits sparse representation if

$$\|\hat{x}\|_0 \ll n.$$

Next we solve the deflated problem

$$\hat{D}\hat{x} = b.$$

# IAS algorithm for sparsity promotion

**Given:** matrix  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ ,  $\beta, \vartheta_j$ ,  $1 \leq j \leq n$ ,

**Initialize:**  $\theta_j^{(0)} = \vartheta_j$ ,  $1 \leq j \leq n$ , set the counter  $k = 0$ .

**Iterate** until stopping criterion is met:

- 1 Update  $x$ : Set  $\theta_j = \theta_j^{(k)}$ ,  $1 \leq j \leq n$  and set

$$x^{(k+1)} = \arg \min \left\{ \|b - Ax\|_2^2 + \|D_{\theta^{(k)}}^{-\frac{1}{2}} x\|_2^2 \right\}.$$

- 2 Update  $\theta_j$ ,  $1 \leq j \leq n$  by solving

$$\frac{\partial \mathcal{E}}{\partial \theta_j}(x^{(k+1)}, \theta) = 0,$$

yielding

$$\theta_j^{(k+1)} = \frac{\vartheta_j}{2} \left( \eta + \sqrt{\eta^2 + \frac{[x_j^{(k+1)}]^2}{4\vartheta_j}} \right), \quad \eta = \beta - \frac{3}{2}.$$

- 3 Advance the counter  $k \rightarrow k + 1$  and check convergence.

# Properties of IAS

When using gamma distribution

- The functional minimized by the IAS algorithm is **convex**.
- The IAS iterates **converge to the unique minimizer**.
- In the limit as  $\eta = \beta - 3/2 \rightarrow 0^+$ , the solution computed with the IAS algorithm **converges to a scaled  $\ell_1$ -regularized solution**,

$$x_{\ell_1} = \operatorname{argmin} \left\{ \frac{1}{2} \|b - Ax\|^2 + \sqrt{2} \sum_{j=1}^n \frac{|x_j|}{\sqrt{\vartheta_j}} \right\},$$

- The value of the parameter  $\beta$  can be selected depending on **how sparse** the solution is believed to be.

# IAS with Group sparsity and Structural Prior for Class Detection

Want to find an approximate representation

$$b = Wh + \varepsilon',$$

where

$$W = [ W^{(1)} \mid \dots \mid W^{(k)} ], \quad h = \begin{bmatrix} h^{(1)} \\ \vdots \\ h^{(k)} \end{bmatrix}, \quad h^{(j)} \in \mathbb{R}^{n_j},$$

with the prior belief that most of the vectors  $h^{(j)}$  are insignificant.

# Structural Prior

Recalling structural covariance matrices  $C^{(j)} \in \mathbb{R}^{n_j \times n_j}$ , we define a group sparsity promoting structural prior for  $h$ ,

$$\begin{aligned} \pi_{H|\Theta}(h | \theta) &\propto \prod_{j=1}^K \left( \frac{1}{\det(\theta_j C^{(j)})} \right)^{1/2} \exp \left( - \sum_{j=1}^K \frac{[h^{(j)}]^\top [C^{(j)}]^{-1} h^{(j)}}{2\theta_j} \right) \\ &\propto \exp \left( - \frac{1}{2} \sum_{j=1}^K n_j \log \theta_j - \frac{1}{2} \sum_{j=1}^K \frac{\|h^{(j)}\|_{C^{(j)}}^2}{\theta_j} \right) \\ &\propto \exp \left( - \frac{1}{2} \sum_{j=1}^K n_j \log \theta_j - \frac{1}{2} \|D_\theta^{-1/2} h\|_2^2 \right), \end{aligned}$$

where  $D_\theta$  is the block diagonal matrix

$$D_\theta = \begin{bmatrix} \theta_1 C_1 & 0 & \dots & 0 \\ 0 & \theta_2 C_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \theta_k C_k \end{bmatrix}.$$

# IAS with Structural Prior and Group Sparsity

- To promote group sparsity, most  $\theta_j$  should be close to zero, with a few large outliers.
- Assume that the variables  $\theta_j$  are independent and distributed according to a gamma distribution.
- Combining the likelihood, enhanced by the dictionary compression error, we obtain a posterior density of the form

$$\begin{aligned} \pi(h, \theta \mid b) \propto \exp \left( -\frac{1}{2} \|C_{\text{DCE}}^{-12} (b - \mu_{\text{DCE}} - Wh)\|^2 - \frac{1}{2} \sum_{j=1}^k \frac{\|h^{(j)}\|_{\ell_2}^2}{\theta_j} \right. \\ \left. + 2 \sum_{j=1}^k \eta_j \log \theta_j - \sum_{j=1}^K \frac{\theta_j}{\vartheta_j} \right), \\ \eta_j = \beta - \frac{n_j + 2}{2}, \quad 1 \leq j \leq K. \end{aligned}$$

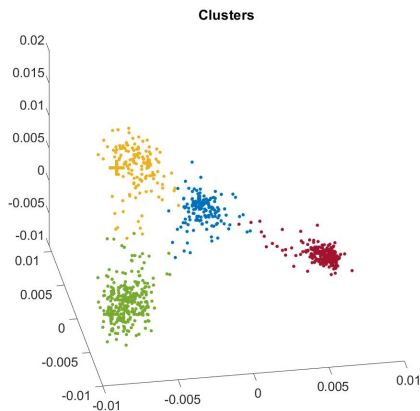
Compute MAP with Group Sparsity IAS (GS-IAS), a version of the IAS algorithm that promotes group sparsity.



# Computed example

The dictionary  $A$  consists of  $n = 834$   $16 \times 16$  digits 0, 1, 3, 9 from the MNIST database, clustered into  $k = 4$  classes using the K-medoid algorithm.

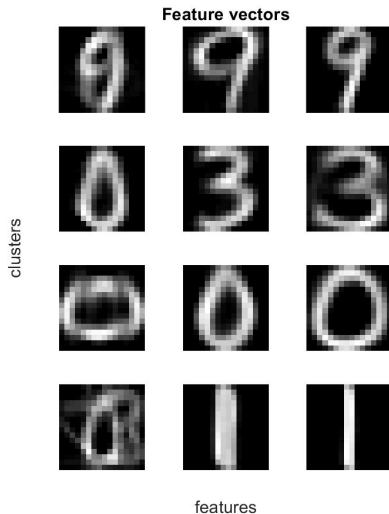
In each test  $b$  is a convex combination of  $1 \leq l \leq 3$  atoms, with at most one atom per cluster.



# Feature extraction and dictionary reduction

Extract  $p = 3$  feature vectors per cluster and assemble the reduced dictionary  $W \in \mathbb{R}^{256 \times 12}$ .

The feature vectors capture the features of two different digits in certain classes where the clustering is not perfect.



# Evaluation: residual and support retrieval

- 4000 tests problems, data convex combinations of 1 2, or 3 digits from different classes.
- Define the thresholded supports of the vectors,

$$S_\delta = \{j \mid x_j > \delta\}, \quad \widehat{S}_\delta = \{j \mid \widehat{x}_j > \delta\}, \quad \delta > 0.$$

- Define the dissimilarity index as the cardinality of the symmetric difference of them,

$$I_\delta(x, \widehat{x}) = \text{card}(S_\delta \Delta \widehat{S}_\delta) = \text{card}((S_\delta \cup \widehat{S}_\delta) \setminus (S_\delta \cap \widehat{S}_\delta)).$$

- When  $I_\delta(x, \widehat{x}) = 0$ , the algorithm has correctly identified the atoms that constitute the data.

	3904 (97.6%)	34 (0.009%)	1	1	59 (0.01%)
$I_\delta(x, \widehat{x})$	0	1	2	3	4

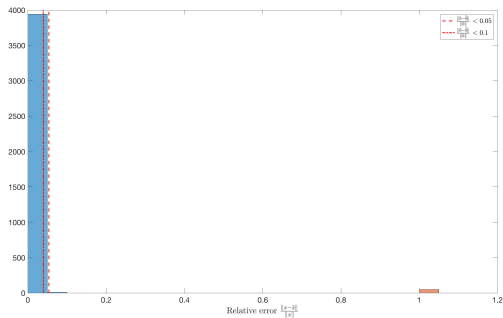
# Fault detection

- Even if some incorrect atoms are identified, the effect on the relative approximation error

$$r(x, \hat{x}) = \frac{\|b - D\hat{x}\|}{\|b\|} = \frac{\|D(x - \hat{x})\|}{\|Dx\|}$$

may be insignificant.

- The number of serious failures of the algorithm can reliably be detected by counting the cases where  $r(x, \hat{x}) > \tau$ , where  $0 < \tau < 1$ .
- The computation of the relative error does not require the knowledge of the true  $x$ .
- The difference between the data  $b$  and the reconstruction  $D\hat{x}$  is a good proxy.
- Method is robust with respect to  $\tau$ : in clearly misidentified cases  $r(x, \hat{x})$  is large.



The mark lines are for  $\tau = 0.05$  (dashed) and  $\tau = 0.1$  (dot-dashed). It is clear that the fault detection mechanism proposed is quite robust with respect to  $\tau$ .

# Dictionary Compression Error

- The data corresponds to a single atom out of a single class.
- Perturb the data by Gaussian white noise with standard deviation 0.03.
- We run the algorithm 5000 times, with and without correcting for DCE.
- Including the DCE leads to a correct cluster identification in 4835 cases, or 97% of the cases.
- Leaving out the DCE, the correct class is identified 4705 times, or 0.94% of the cases.
- Once the correct class is identified, the correct atom is identified in both protocols in approximately 98% of the cases.
- The effect of the DCE is observable but not as crucial as in the next example.

# LIGO/VIRGO glitches classification

- Observed gravitational wave, or an event of any unknown origin registered by the device, is to be identified based on a precomputed dictionary, as finding a universal predictive parametric model is not feasible.
- Data stream is contaminated by frequent artifacts of non-gravitational origin, known as **glitches**, that due to their high occurrence rate may lead to false coincidence detections.
- The glitches are transients with particular morphologies that allow them to be classified and identified by using precomputed glitch libraries.

# Glitches dictionary

Consider three classes of glitches, sine Gaussian (SG), Gaussian (G) and ring down (RD), parametrized as

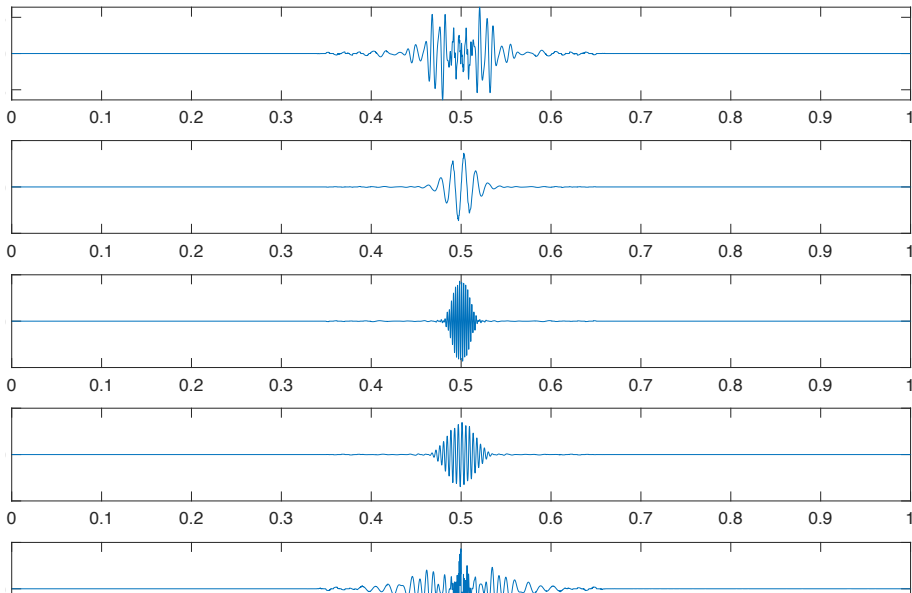
$$\begin{aligned} h_{\text{SG}}(t) &= h_0 \sin(2\pi f_0(t - t_0)) e^{-(t-t_0)^2/2\tau^2}, \\ h_{\text{G}}(t) &= h_0 e^{-(t-t_0)^2/2\tau^2}, \\ h_{\text{RD}}(t) &= h_0 \sin(2\pi f_0(t - t_0)) e^{-(t-t_0)/\tau} \theta(t - t_0), \end{aligned}$$

where  $h_0$  is the amplitude,  $f_0$  is the frequency,  $t_0$  is the center time and  $\tau$  is the characteristic time of the glitch, and  $\theta$  is the Heaviside function.

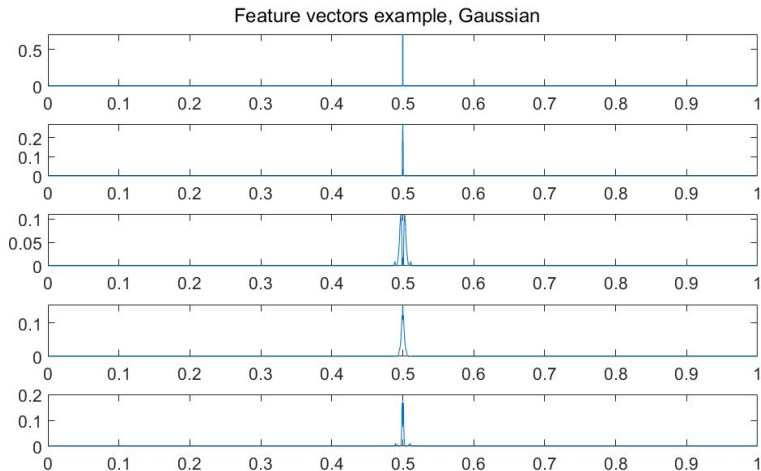
- The dictionary is made of 1000 random glitches per class.
- Dimensionality of atoms is significant
- Number of feature vectors changes for the three classes : 50 for sine Gaussian and ring down signals, 11 for Gaussian signals.



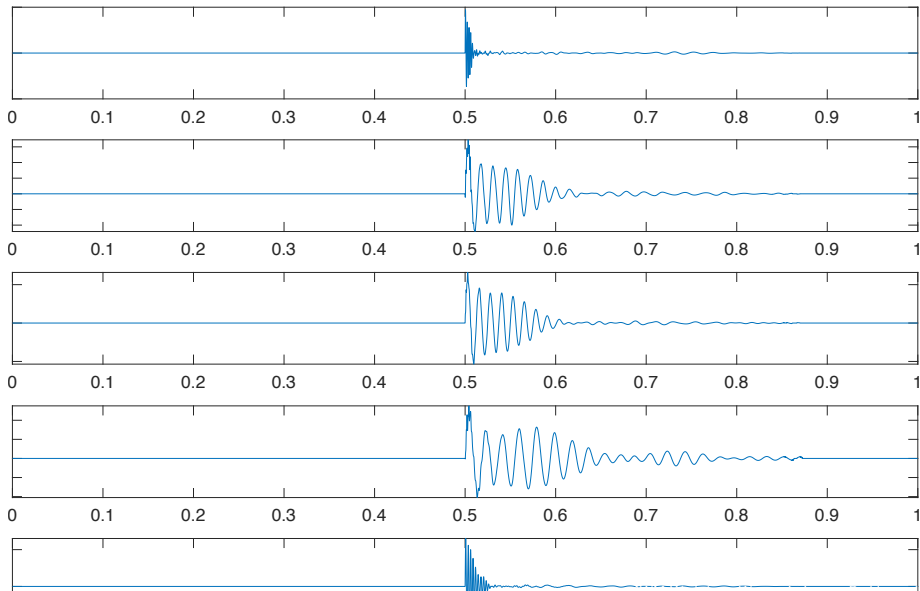
# Sine Gaussians feature vectors



# Gaussians feature vectors



# Ring down feature vectors



To identify the class, for a given data vector  $b$ , we find the coefficient vector  $h^{(j)}$  for each cluster,  $1 \leq j \leq 3$ , such that

$$b \approx W^{(j)} h^{(j)} = \widehat{b}^{(j)},$$

where  $h^{(j)}$  is computed using the IAS algorithm with sparsity and positivity constraints, including the DCE correction. Then compare  $\widehat{b}^{(j)}$  with  $b$  using the structural dissimilarity index,

$$\text{DSSIM}(b, \widehat{b}^{(j)}) = \frac{1}{2}(1 - \text{SSIM}(b, \widehat{b}^{(j)})),$$

where SSIM is the structural similarity index. The chosen class is the one corresponding to smallest dissimilarity index.

# Results

We run the test 20000 times with additive noise  $\sigma = 0.02$ , and find that the algorithm identifies correctly the class with 100% accuracy (left). Increasing the additive noise level to  $\sigma = 0.05$  still get 99.8% correct classification (right).

	SG	G	RD
SG	6812	0	0
G	0	6546	0
RD	0	0	6642

	SG	G	RD
SG	6798	0	0
G	0	6522	37
RD	1	0	6642

# Importance of DCE

For comparison, we then run the same experiment leaving out the DCE correction. Even with the lower noise level,  $\sigma = 0.02$ , the class identification goes egregiously astray, giving a 100% misclassification result, with the confusion matrix

	SG	G	RD
SG	0	6754	0
G	6592	0	0
RD	6654	0	0

# Simplified DCE

Finally, we approximate the DCE covariance matrix by a scaled identity matrix,

$$\mathbf{C}_{\text{DCE}}^{(j)} \approx \text{trace}(\mathbf{C}_{\text{DCE}}^{(j)}) \mathbf{I}_m.$$

With the noise level  $\sigma = 0.02$ , the identification of the clusters was 100% correct, and with  $\sigma = 0.05$ , an identification performance of 99.5% was achieved again.

# References

- Llorens-Monteagudo M, Torres-Forné A, Font JA and Marquina A (2019) Classification of gravitational-wave glitches via dictionary learning. *Classical and Quantum Gravity* **36** 075005
- Pragliola M, Calvetti D and Somersalo E (2022) Overcomplete representation in a hierarchical Bayesian framework. *Inverse Problems and Imaging* **16**: 19-38. doi: 10.3934/ipi.2021039
- Waniorek N, Somersalo E and Calvetti D (2023) Bayesian hierarchical dictionary learning. *Inverse Problems* (electronic version: DOI 10.1088/1361-6420/acad21)

The details of this presentation are in a forthcoming manuscript

- A. Bocchinfuso, D. Calvetti and E. Somersalo (2023) Error enhanced layered Bayesian dictionary learning.