# Inverse problems and optimization

#### Nelly PUSTELNIK CNRS, Laboratoire de Physique de l'ENS Lyon Laurent CONDAT CNRS, Gipsa-lab

Course 4: Sparsity

January, 19th 2017

Frame	ℓ <sub>1</sub> -norm	CS	$\substack{\ell_{1,p}\text{-norm}\\00000000}$
•000000	0000000000	00000000	
			2/36

# Motivation: image restoration





Frame	ℓ <sub>1</sub> -norm	CS	$\ell_{1,p}$ -norm
•000000	0000000000	00000000	
			2/36

# Motivation: image restoration







Frame
000000

CS 000000000  $\ell_{1,p}$ -norm

3/36

# Sparse transforms: finite difference

![](_page_3_Figure_6.jpeg)

![](_page_3_Figure_7.jpeg)

Frame	$\ell_1$ -norm	CS	$\ell_{1,p}$ -norm
000000	0000000000	00000000	000000
			4/36

# Sparse transforms: finite difference

![](_page_4_Figure_2.jpeg)

- ▶ Horizontal and vertical difference filters:  $f_1 = [1 1]$  and  $f_2 = f_1 \top$ ,
- K = 2N,
- Link between  $(f_1, f_2)$  and  $(F_1, F_2)$ : cf. slides10-12 course 1,
- ▶ Resulting  $F = [F_1^\top, F_2^\top]^\top$  and  $F^* = F^\top$ ,
- ►  $F^*F \neq \text{Id.}$

CS 000000000  $\ell_{1,p}$ -norm 00000000 5/36

# Sparse transforms: wavelet transform

![](_page_5_Figure_5.jpeg)

![](_page_5_Figure_6.jpeg)

CS 000000000 ℓ<sub>1,p</sub>-norm 00000000 6/36

# Sparse transforms: wavelet transform

![](_page_6_Figure_5.jpeg)

 F (resp. F\*): Concatenation of matrices associated to filtering and decimation (resp. upsampling) operations.

![](_page_6_Figure_7.jpeg)

Frame
0000000

CS 000000000  $\ell_{1,p}$ -norm

7/36

# Sparse transforms: frame transform

![](_page_7_Figure_6.jpeg)

![](_page_7_Figure_7.jpeg)

Frame	$\ell_1$ -norm	CS	$\ell_{1,p}$ -norm
000000	0000000000	00000000	0000000
			8/36

# Sparse transforms: dictionary

 Dictionary: set of elementary signals

 $D = [d_1, \ldots, d_K] \in \mathbb{R}^{N \times K}$ 

- D is "adapted" to x if it can represent it with a few elements, i.e., there exists a sparse vector θ in ℝ<sup>K</sup> such that x ~ Dθ.
- $\triangleright$   $\theta$ : sparse code.

![](_page_8_Figure_6.jpeg)

![](_page_8_Figure_7.jpeg)

Frame	$\ell_1$ -norm	CS	$\ell_{1,p}$ -norm
0000000	●000000000	00000000	00000000
			9/36

# Sparse transforms

**Quadratic regularization** - Wiener filtering [Wiener,1949]

$$\widehat{y} = \operatorname*{argmin}_{y \in \mathbb{R}^N} \|Hx - z\|_2^2 + \lambda \|Fx\|_2^2 \qquad \lambda > 0$$

![](_page_9_Picture_4.jpeg)

Use of frames - Soft-thresholding  
[Haar,1910] [Mallat,2009]  

$$\widehat{y} = \underset{y \in \mathbb{R}^{N}}{\operatorname{argmin}} \|Hx - z\|_{2}^{2} + \lambda \frac{\|Fx\|_{1}}{\|Fx\|_{1}} \quad \lambda > 0$$

![](_page_9_Picture_6.jpeg)

► The squared ℓ<sub>2</sub>-norm induces "smoothness" while ℓ<sub>1</sub>-norm induces sparsity. [Chen et al., 1999, Tibshirani, 1996]

Frame	$\ell_1$ -norm	CS	$\ell_{1,p}$ -norm
			9/36

# Sparse transforms

**Quadratic regularization** - Wiener filtering [Wiener, 1949]

$$\widehat{y} = \operatorname*{argmin}_{y \in \mathbb{R}^N} \|Hx - z\|_2^2 + \lambda \|Fx\|_2^2 \qquad \lambda > 0$$

 $\widehat{y} = \operatorname{argmin} \| Hx - z \|_2^2 + \lambda \| Fx \|_1$   $\lambda > 0$ 

**Use of frames** - Soft-thresholding

[Haar,1910] [Mallat,2009]

 $v \in \mathbb{R}^N$ 

![](_page_10_Picture_4.jpeg)

![](_page_10_Picture_5.jpeg)

► The squared ℓ<sub>2</sub>-norm induces "smoothness" while ℓ<sub>1</sub>-norm induces sparsity. [Chen et al., 1999, Tibshirani, 1996]

Frame	ℓ <sub>1</sub> -norm	CS	$\ell_{1,p}$ -norm
0000000	0●000000000	00000000	
			10/36

# Sparse transforms: analysis/synthesis

# Analysis formulation

$$\widehat{x} \in \underset{x \in \mathbb{R}^{N}}{\operatorname{Argmin}} \|Hx - z\|_{2}^{2} + \lambda \|Fx\|_{1} \qquad \lambda > 0$$

#### Synthesis formulation

$$\widehat{x} = F^* \widehat{\theta} \quad \text{with} \quad \widehat{\theta} \in \operatorname*{Argmin}_{\theta \in \mathbb{R}^K} \|HF^* \theta - z\|_2^2 + \lambda \|\theta\|_1 \qquad \lambda > 0$$

Sparse coding with  $F^* = D$ .

 Equivalence for some specific conditions over F such orthogonality. [Pustelnik et al. 2016]

Frame 0000000	ℓ <sub>1</sub> -norm 00●0000000	CS 00000000	$\ell_{1,p}$ -norm 00000000
			11/36

# Variational approach: Bayesian formulation

• 
$$u = Hx = (u^{(i)})_{1 \le i \le N}$$
: realization of a random vector  $U$ .

- z: realization of a random vector Z.
- $\theta = Fx = (\theta^{(k)})_{1 \le k \le K}$ : realization of a random vector  $\overline{\Theta} = (\Theta^{(k)})_{1 < k < K}$  having independent components.

# MAP estimator (Maximum A Posteriori) $\max_{x} P(U = Hx | Z = z)$ $\max_{\theta} P(Z = z | U = HF^{*}\theta) \cdot P(\Theta = \theta)$ $\min_{\theta} -\underbrace{\ln P(Z = z | U = HF^{*}\theta)}_{\text{Data fidelity}} -\underbrace{\sum_{k=1}^{K} \ln p_{\Theta^{(k)}}(\theta^{(k)})}_{\text{A priori}}$

Frame	$\ell_1$ -norm	CS	$\ell_{1,p}$ -norm
000000	0000000000	00000000	0000000

12/36

# Variational approach: Bayesian formulation

![](_page_13_Figure_3.jpeg)

![](_page_14_Figure_0.jpeg)

where

$$\mathsf{P}(Z = z \mid U = AF^*\theta) = \frac{1}{(2\pi\alpha)^{M/2}} \exp\left\{-\frac{\|HF^*\theta - z\|_2^2}{2\alpha}\right\}$$

and

$$p_{\Theta^{(k)}}(\theta^{(k)}) = \frac{1}{C_k} \exp\{-\lambda_k |\theta^{(k)}|\}$$

$$\min_{\theta} \quad \frac{1}{2\alpha} \|HF^*\theta - z\|_2^2 + \sum_{k=1}^K \lambda_k |\theta^{(k)}|$$

Frame 0000000	ℓ <sub>1</sub> -norm 00000●00000	CS 00000000	$\ell_{1,p}$ -norm
			14/36

# Variational approach: Equivalent formulations

- Geophysics: [Claerbout and Muir, 1973, Taylor et al., 1979],
- Statistics: Lasso [Tibshirani, 1996]
- Signal processing: Basis pursuit [Chen et al., 1999]

Equivalent formulations:

$$\begin{split} \min_{\theta \in \mathbb{R}^{K}} & \frac{1}{2\alpha} \| HF^{*}\theta - z \|_{2}^{2} + \lambda \| \theta \|_{1} \\ \min_{\theta \in \mathbb{R}^{K}} & \frac{1}{2\alpha} \| HF^{*}\theta - z \|_{2}^{2} \quad \text{s.t.} \quad \| \theta \|_{1} \leq \mu \\ \min_{\theta \in \mathbb{R}^{K}} & \| \theta \|_{1} \quad \text{s.t.} \quad \| HF^{*}\theta - z \|_{2}^{2} \leq \varepsilon \end{split}$$

Frame	$\ell_1$ -norm	CS	$\ell_{1,p}$ -norm
0000000	0000000000	00000000	0000000
			15/36

![](_page_16_Picture_2.jpeg)

(figures extracted from J. Mairal course)

Frame 0000000	ℓ <sub>1</sub> -norm 0000000●000	<b>CS</b> 00000000	$\ell_{1,p}$ -norm
			16/36

$$\min_{\theta \in \mathbb{R}} \ \frac{1}{2} (\theta - z)^2 + \lambda |\theta|$$

- Piecewise quadratic function with a non-differentiability at zero.
- Fermat rule for non-differentiable function:

$$0 \in \widehat{\theta} - z + \lambda u$$
 with  $u \in \partial |\widehat{\theta}|$ 

This lead to soft-thresholding

$$\widehat{\theta} = \operatorname{sign}(z) \max(0, |z| - \lambda)$$

Frame 0000000	ℓ <sub>1</sub> -norm 00000000●00	CS 00000000	$\ell_{1,p}$ -norm
			17/36

► Resolution of  $\min_{x \in \mathbb{R}^{K}} ||x||_{1}$  s.t. Ax = y

![](_page_18_Figure_3.jpeg)

Frame	$\ell_1$ -norm	CS	$\ell_{1,p}$ -norm
0000000	0000000000	00000000	0000000
			18/36

# ▶ $\ell_2^2$ versus $\ell_1$

![](_page_19_Figure_3.jpeg)

![](_page_20_Figure_0.jpeg)

# Non-convex functions

▶  $\ell_q$ -norms with  $q \in ]0, 1[$  [Frank and Friedman, 1993]

$$\psi(\theta) = \sum_{k=1}^{K} |\theta_k|^q$$

Log penalty [Candès et al. 2008]

$$\psi(\theta) = \sum_{k=1}^{K} \log(|\theta_k| + \varepsilon)$$

- Several others [Nikolova, 2007]
- Non-convex penalties leading to convex criterion [Selesnick et al. 2015]

Frame 0000000	$\ell_1$ -norm	CS ●00000000	$\ell_{1,p}$ -norm 00000000
			20/36

## Compressed sensing

#### Sparsity assumption

$$heta_{0} \in \operatorname*{arg\,min}_{ heta \in \mathbb{R}^{K}} rac{1}{2lpha} \|D heta - z\|_{2}^{2} + \lambda \| heta\|_{0}$$

where  $(\forall \theta = (\theta_k)_{1 \le k \le K} \in \mathbb{R}^K) \|\theta\|_0 \equiv \#\{k : \theta_k \neq 0\}$ 

#### **Convex relaxation**

$$\widehat{ heta}_0 \in \operatorname*{arg\,min}_{ heta \in \mathbb{R}^K} rac{1}{2lpha} \|D heta - z\|_2^2 + \lambda \| heta\|_1$$

where  $(\forall \theta = (\theta_k)_{1 \le k \le K} \in \mathbb{R}^K) \|\theta\|_1 = \sum_{k=1}^N |\theta_k|$ 

#### ⇒ Compressed sensing [Donoho,2004][Candès et al. 2006]

![](_page_22_Figure_0.jpeg)

# Compressed sensing

#### Restricted Isometry Property (RIP) [Candès, 2008]:

$$(1 - \delta_{2s}^{\min}) \|\theta\|^2 \le \|D\theta\|^2 \le (1 + \delta_{2s}^{\max}) \|\theta\|^2$$

where

$$\begin{cases} \delta_s^{\min} = 1 - \lambda^{\min}(D_l^* D_l) \\ \delta_s^{\max} = \lambda^{\max}(D_l^* D_l) - 1 \end{cases} \quad \text{and} \quad D_l = (d_i)_{i \in I} \text{ with } |I| = s. \end{cases}$$

#### Theorem [Candès, 2008]

Assume that  $\delta_{2s} = \min(\delta_{2s}^{\min}, \delta_{2s}^{\max}) < \sqrt{2} - 1$ . Then the solution  $\hat{x}$  to the  $\ell_1$ -minimization obeys

 $\|\hat{\theta} - \theta_0\|_1 < C_0 \|\theta_0 - \theta_s\|_1$  and  $\|\hat{\theta} - \theta_0\|_2 \le C_0 s^{-1/2} \|\theta_0 - \theta_s\|_1$ for some constant  $C_0$  and where  $\theta_s$  denotes the s highest non-zero component of  $\theta_0$ . In particular, if  $\theta$  is s-sparse, the recovery is exact.

Frame 0000000	ℓ <sub>1</sub> -norm 0000000000	CS ००●००००००	$\ell_{1,p}$ -norm 00000000
			22/36
Compressed sen	sing		

- Restricted Isometry Property (RIP) [Candès, 2008]:
  - Extended proposition to insure robustness to noise.
  - Sufficient condition (not necessary).
  - NP-hard to compute  $\delta_{2s}$ .

•  $\delta_{2s} < \sqrt{2} - 1$ : restrictive condition even for random matrices [Dossal et al., 2010].

Frame 0000000	ℓ <sub>1</sub> -norm 0000000000	CS ०००●०००००	$\ell_{1,p}$ -norm
			23/36
<b>c</b>			

# Compressed sensing

# Sub-differential of l<sub>1</sub> norm [Fuchs, 2004][Tropp, 2006]

#### Theorem [Fuchs, 2004]

The solution  $\theta_0$  of  $D\theta = z$  can be recovered from the unique optimum point  $\hat{\theta}$  of  $\ell_1$ -minimization if

1)  $|d_j^{\top} d_0| < 1$ ,  $\forall d_j \notin \overline{D}_0$  with  $d_0 = \overline{D}_0^{+\top} \operatorname{sign}(\overline{\theta}_0)$ , 2)  $h \in ]0, h_m[$  the domain in which  $\operatorname{sign}\{\overline{\theta}_0 - h(\overline{D}_0^{\top}\overline{D}_0)^{-1}\operatorname{sign}(\overline{\theta}_0)\} = \operatorname{sign}(\overline{\theta}_0),$ 

where

- $\overline{\theta}_0$  is built upon the nonzero component of  $x_0$ ,
- $\overline{D}_0$  a full rank matrix such that  $Dx_0 = \overline{D}_0 \overline{x}_0$ ,
- $\overline{D}_0^+ = (\overline{D}_0^\top \overline{D}_0)^{-1} \overline{D}_0^\top$  denotes the pseudo-inverse of  $\overline{D}$ .

Frame 0000000	ℓ <sub>1</sub> -norm 0000000000	CS 0000●0000	$\ell_{1,p}$ -norm
			24/36
Compressed sen	sing		

- ▶ Sub-differential of ℓ<sub>1</sub> norm [Fuchs, 2004][Tropp, 2006]
  - Extended proposition to ensure robustness to noise.
  - Sufficient condition (not necessary).
  - NP-hard to check the condition  $|d_j^{\top} d_0| < 1$ ,  $\forall d_j \notin \overline{D}_0$  with  $d_0 = \overline{D}_0^{+\top} \operatorname{sign}(\overline{\theta}_0)$  for all vector *s*-sparse.

Frame 0000000	ℓ <sub>1</sub> -norm 0000000000	CS 00000●000	$\ell_{1,p}$ -norm
			25/36

# Compressed sensing

Coherence [Fuchs, 2004]

#### Theorem [Fuchs, 2004]

The solution  $\theta_0$  of  $D\theta = z$  can be recovered from the unique optimum point  $\hat{\theta}$  of  $\ell_1$ -minimization if

$$\| heta_0\|_0 < rac{1}{2}\Big(1+rac{1}{M}\Big) \qquad ext{where} \qquad M = \max_{1 \leq i \neq j \leq m} |d_i^{ op} d_j|.$$

Frame 0000000	ℓ <sub>1</sub> -norm 0000000000	CS ○○○○○○●○○	$\ell_{1,p}$ -norm 00000000
			26/36
Compressed sen	sing		

- Coherence [Fuchs, 2004]
  - Extended proposition to ensure robustness to noise.
  - Sub-differential of  $\ell_1$  norm condition  $\Rightarrow$  Coherence condition.
  - Sufficient condition (not necessary).

• Easy to compute  $\frac{1}{2}\left(1+\frac{1}{M}\right)$  and thus to know the allowed sparsity degree.

![](_page_28_Figure_0.jpeg)

#### Polytope theory [Donoho,2004] : NSC

- A large panel of the literature provides results based on greedy algorithms (MP, OMP, OLS) solving rather that I1-minimization. MP [Mallat 93], projection pursuit [Friedman 81], [Huber 85], pure greedy [Temlyakov 08], OMP : [Pati 93], [Zhang 93], [Davis 94], orthogonal greedy algorithm [Temlyakov 08], OLS : [Chen 89], forward selection [Miller 02], greedy algorithm [Natarajan 95], order recursive matching pursuit [Cotter 99], optimized orthogonal matching pursuit [Reibollo-Neira 02], pure orthogonal matching pursuit [Foucart 11].
- ► Spectral sparsity [Fazel et al., 2001] [Srebro et al., 2005] : ℓ<sub>1</sub>-ℓ<sub>0</sub> formulation for matrices. Minimization of the rank or its relaxation involving nuclear norm

![](_page_29_Figure_0.jpeg)

#### Difference between random and deterministic matrices:

- In a random context, the eigenvalue distribution can be controlled and it results an explicit relation between *s*, *K*, and *N*.
- In a deterministic context and without prior onto the support or the sign, the validation of the previous properties (RIP, Coherence,...) requires a high computational cost.

Frame 0000000	ℓ <sub>1</sub> -norm 0000000000	CS 00000000	$\ell_{1,p}$ -norm
			29/36
<u> </u>			

# Structured sparsity

![](_page_30_Picture_2.jpeg)

(figures extracted from J. Mairal course)

![](_page_31_Figure_0.jpeg)

#### Anisotropic total variation [Rudin et al. 1992]

$$\widehat{x} = \underset{x \in \mathbb{R}^{N}}{\arg\min} \|x - z\|_{2}^{2} + \lambda \sum_{n=1}^{N-1} |x_{n+1} - x_{n}|$$

- Statistics: fused Lasso
   [Tibshirani et al 2005]
- *l*<sub>0</sub>-formulation: Potts model [Geman, Geman 1984]
   [Yao 1984]

[Mumford, Shah 1989]

[Winkler, Liebscher 2002]

![](_page_31_Figure_8.jpeg)

(cf. [Sowa et al. 2005])

![](_page_32_Figure_0.jpeg)

# Anisotropic total variation [Rudin et al. 1992] $\widehat{x} = \underset{x \in \mathbb{R}^{N}}{\arg \min} \|x - z\|_{2}^{2} + \lambda \sum_{n_{1}=1}^{N-1} \sum_{n_{2}=1}^{N-1} (|x_{n_{1}+1,n_{2}} - x_{n_{1},n_{2}}|^{2} + |x_{n_{1},n_{2}+1} - x_{n_{1},n_{2}}|^{2})$

- ▶ Horizontal and vertical difference filters:  $f_1 = [1 - 1]$  and  $f_2 = f_1 \top$
- Link between (f<sub>1</sub>, f<sub>2</sub>) and (F<sub>1</sub>, F<sub>2</sub>): cf. slides 10-12 course 1
- ► Sparse transform:  $F = [F_1^\top, F_2^\top]^\top \in \mathbb{R}^{2N \times N}$

• Regularization:  $\psi(x) = \|Fx\|_1 = \|F_1x\|_1 + \|F_2x\|_1$ 

![](_page_32_Figure_7.jpeg)

![](_page_33_Figure_0.jpeg)

Isotropic total variation [Rudin et al. 1992]  

$$\widehat{x} = \underset{x \in \mathbb{R}^{N}}{\arg \min} \|x - z\|_{2}^{2} + \lambda \sum_{n_{1}=1}^{N-1} \sum_{n_{2}=1}^{N-1} \sqrt{|x_{n_{1}+1,n_{2}} - x_{n_{1},n_{2}}|^{2} + |x_{n_{1},n_{2}+1} - x_{n_{1},n_{2}}|^{2}}$$

Horizontal and vertical difference filters:

 $\mathit{f}_1 = [1-1]$  and  $\mathit{f}_2 = \mathit{f}_1 op$ 

• Link between  $(f_1, f_2)$  and  $(F_1, F_2)$ :

cf. slides 10-12 course 1

► Sparse transform:  $F = [F_1^\top, F_2^\top]^\top \in \mathbb{R}^{2N \times N}$ 

• Regularization:  $\psi(x) = ||Fx||_{1,2}$ 

![](_page_33_Figure_9.jpeg)

Frame 0000000	ℓ <sub>1</sub> -norm 0000000000	CS 00000000	$\ell_{1,p}$ -norm 0000 $\bullet$ 000
			33/36

![](_page_34_Picture_2.jpeg)

![](_page_34_Picture_3.jpeg)

Anisotropic TV

![](_page_34_Picture_5.jpeg)

Isotropic TV

![](_page_34_Picture_7.jpeg)

Frame 0000000	$\ell_1$ -norm	<b>CS</b> 00000000	$\ell_{1,p}$ -norm
			33/36

![](_page_35_Figure_2.jpeg)

Noisy z

![](_page_35_Picture_4.jpeg)

Anisotropic TV

![](_page_35_Picture_6.jpeg)

![](_page_35_Figure_7.jpeg)

![](_page_35_Picture_8.jpeg)

![](_page_36_Figure_0.jpeg)

# Mixed-norm

**Group-lasso** [Turlach et al., 2005] [Yuan and Lin, 2006] [Zhao et al., 2009] [Grandvalet and Canu, 1999] [Bakin, 1999]

$$\widehat{x} = \underset{\theta \in \mathbb{R}^{K}}{\arg\min} \|D\theta - z\|_{2}^{2} + \lambda \sum_{g \in \mathcal{G}} \|\theta_{g}\|_{q}$$

- $\mathcal{G}$  is a partition of  $\{1, \ldots, K\}$
- ▶  $q = \{2, +\infty\}$
- ► Can be interpreted as the ℓ<sub>1</sub>-norm of (||θ<sub>g</sub>||<sub>q</sub>)<sub>g∈G</sub>.

 $\rightarrow$  Non-overlapping groups

![](_page_36_Figure_8.jpeg)

Frame 0000000	$\ell_1$ -norm	CS 00000000	$\ell_{1,p}$ -norm
			35/36
Structured	sparsity		

- Tree-structure [Zhao et al. 2009]
- Select a union of groups [Jacob et al. 2009]
- Zero-pattern in a union of groups [Jenatton et al, 2011]
  - $\rightarrow$  Overlapping groups

![](_page_37_Figure_5.jpeg)

![](_page_38_Figure_0.jpeg)

- Sparsity is not always good. If possible, try  $\ell_2^2$  before trying  $\ell_1$ .
- The dictionaries used in practice rarely satisfy the assumptions ensuring sparse recovery.
- Dictionary learning consists of estimating D and  $\theta$  simultaneously.
- There are numerous ways of designing sparse regularization functions adapted to a particular problem. Choosing the best one is not easy and requires some domain knowledge.
- Solving criterion involving l<sub>1</sub>-norm requires specific algorithmic procedures described in next courses.