Structured Sparsity-Inducing Norms: Statistical and Algorithmic Properties with Applications to Neuroimaging

Rodolphe Jenatton

SIERRA project, INRIA Rocquencourt - Ecole Normale Supérieure

PhD defense at ENS Cachan, November 24, 2011

Advisors: Jean-Yves Audibert Francis Bach Reviewers: Laurent El Ghaoui Massimiliano Pontil



Examinateurs: Rémi Gribonval Eric Moulines Guillaume Obozinski Bertrand Thirion

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

Acknowledgments









Francis Bach

Jean-Yves Audibert Guillaume Obozinski Julien Mairal









Bertrand Thirion Alexandre Gramfort Gaël Varoqueaux



Rémi Gribonval

Internetions pretametions

SIERRA/WILLOW

Vincent Michel

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Sparsity

- Important concept in statistics, machine learning,...
 - Favour "simple models"
 - Easier interpretation, cheaper post-processing
 - Models with few parameters, feature selection

Sparsity

- Important concept in statistics, machine learning,...
 - Favour "simple models"
 - Easier interpretation, cheaper post-processing
 - Models with few parameters, feature selection
- Applications:
 - Compressed sensing [Candes and Tao, 2005]
 - Graphical models [Meinshausen and Bühlmann, 2006]
 - Signal/image processing tasks
 - E.g., denoising [Chen et al., 1998; Mairal, 2010]

Sparsity

- Important concept in statistics, machine learning,...
 - Favour "simple models"
 - Easier interpretation, cheaper post-processing
 - Models with few parameters, feature selection
- Applications:
 - Compressed sensing [Candes and Tao, 2005]
 - Graphical models [Meinshausen and Bühlmann, 2006]
 - Signal/image processing tasks
 - E.g., denoising [Chen et al., 1998; Mairal, 2010]
- Disregards structure:
 - Prior knowledge only about cardinality

Structure? The example of neuroimaging

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- In practice, know more than just cardinality!
- Sparsity comes with structural information

Structure? The example of neuroimaging

- In practice, know more than just cardinality!
- Sparsity comes with structural information
- Spatial structure in neuroimaging [Chklovskii and Koulakov, 2004; Gramfort et al., 2011]
 - Sparsity: few discriminative voxels
 - Spatiality: clusters according to the geometry of the brain



• Spatial:

- Bioinformatics: contiguity due to the genome organization [Rapaport et al., 2008]
- Image segmentation: neighboring pixels [Boykov et al., 2001]

- Spatial:
 - Bioinformatics: contiguity due to the genome organization [Rapaport et al., 2008]
 - Image segmentation: neighboring pixels [Boykov et al., 2001]

- Temporal:
 - Time series of gene expressions [Tibau Puig et al., 2011]

- Spatial:
 - Bioinformatics: contiguity due to the genome organization [Rapaport et al., 2008]
 - Image segmentation: neighboring pixels [Boykov et al., 2001]
- Temporal:
 - Time series of gene expressions [Tibau Puig et al., 2011]
- Hierarchical:
 - NLP: hierarchical selection of templates [Martins et al., 2011]
 - Graphical models: potential orders of interactions [Schmidt and Murphy, 2010]

• Wavelet: zero-tree coding [Shapiro, 1993]

- Spatial:
 - Bioinformatics: contiguity due to the genome organization [Rapaport et al., 2008]
 - Image segmentation: neighboring pixels [Boykov et al., 2001]
- Temporal:
 - Time series of gene expressions [Tibau Puig et al., 2011]
- Hierarchical:
 - NLP: hierarchical selection of templates [Martins et al., 2011]
 - Graphical models: potential orders of interactions [Schmidt and Murphy, 2010]

- Wavelet: zero-tree coding [Shapiro, 1993]
- Physics of the sensors:
 - Multipath radar signals [Rebafka et al., 2011]

Some questions addressed by the thesis

- Within a convex framework...
- How to take into account structure while inducing sparsity?
- Which structures can be considered?
- Statistical properties of the resulting estimators?
- Can we design efficient algorithms for practical problems?

• Do we observe improvements in applications?

Topic of Part I

- New family of structured sparsity-inducing norms
- Analysis of
 - Type of structures which can be encoded
 - Automatic norm design
 - Statistical property: low- and high-dimensional consistency

Related publication:

<u>R. Jenatton</u>, J.-Y. Audibert, F. Bach. Structured Variable Selection with Sparsity-Inducing Norms. In *Journal of Machine Learning Research*, 12, 2777-2824. 2011

Topic of Part II

- Efficient optimization methods for structured sparsity:
 - Active-set algorithms
 - Reweighted- ℓ_2 schemes
 - Proximal-gradient techniques
- Efficient computations of proximal operator/dual norms
- Connections with operations research and network flows

Related publications:

R. Jenatton, J.-Y. Audibert, F. Bach. Structured Variable Selection with Sparsity-Inducing Norms. In Journal of Machine Learning Research, 12, 2777-2824. 2011 R. Jenatton, G. Obozinski, F. Bach. Structured sparse principal component analysis. In International Conference

<u>R. Jenatton</u>, G. Obozinski, F. Bach. Structured sparse principal component analysis. In *International Conference* on Artificial Intelligence and Statistics (AISTATS). 2010

<u>R. Jenatton*</u>, J. Mairal*, G. Obozinski, F. Bach. Proximal Methods for Sparse Hierarchical Dictionary Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*. 2010

<u>R. Jenatton</u>^{*}, J. Mairal^{*}, G. Obozinski, F. Bach. Proximal Methods for Hierarchical Sparse Coding. In *Journal of Machine Learning Research*, 12, 2297-2334. 2011

J. Mairal*, <u>R. Jenatton*</u>, G. Obozinski, F. Bach. Network Flow Algorithms for Structured Sparsity. Advances in Neural Information Processing Systems. 2010

J. Mairal*, <u>R. Jenatton</u>*, G. Obozinski, F. Bach. Convex and Network Flow Optimization for Structured Sparsity. In *Journal of Machine Learning Research*, 12, 2681-2720. 2011 (*equal contributions)

Topic of Part III

- Sparse structured dictionary learning/matrix factorization:
 - Exploit previous norms for latent models
 - · Learn representations adapted to structured signals
 - Flexible and rich framework
- Connections with probabilistic topic models [Blei et al., 2003]
- Applications in text/image processing and computer-vision

Related publications:

R. Jenatton, G. Obozinski, F. Bach. Structured sparse principal component analysis. In International Conference on Artificial Intelligence and Statistics (AISTATS). 2010 <u>R. Jenatton*</u>, J. Mairal*, G. Obozinski, F. Bach. Proximal Methods for Sparse Hierarchical Dictionary Learning. In Proceedings of the International Conference on Machine Learning (ICML). 2010 <u>P. Jenetres*</u>, J. Mairal*, C. Obozinski, E. Bach. Deroximal Methods for Microsoften Conference on Jacobian Methods for Sparse Hierarchical Science of the International Conference on Machine Learning (ICML). 2010

<u>R. Jenatton</u>*, J. Mairal*, G. Obozinski, F. Bach. Proximal Methods for Hierarchical Sparse Coding. In *Journal of Machine Learning Research*, 12, 2297-2334. 2011

Topic of Part IV

- Applications of structured sparsity to neuroimaging
- Supervised: prediction of object sizes from fMRI signals
 - Large scale problem
 - Hierarchical norm to have multiscale representations of voxels
 - Gain in robustness for inter-subject validation
- Unsupervised: resting-state brain activity modeling
 - Sparse dictionary learning adapted to 3D structure
 - Generative model for model selection and evaluation

Related publications:

<u>R. Jenatton</u>, A. Gramfort, V. Michel, G. Obozinski, F. Bach, and B. Thirion. Multi-scale Mining of fMRI Data with Hierarchical Structured Sparsity. In International Workshop on Pattern Recognition in Neuroimaging (PRNI). 2011 <u>R. Jenatton</u>, A. Gramfort, V. Michel, G. Obozinski, E. Eger, F. Bach, and B. Thirion. Multi-scale Mining of fMRI Data with Hierarchical Structured Sparsity. Preprint arXiv:1105.0363 Submitted to SIAM Journal on Imaging Sciences. 2011

G. Varoquaux, R. Jenatton, A. Gramfort, G. Obozinski, F. Bach, and B. Thirion. Sparse Structured Dictionary Learning for Brain Resting-State Activity Modeling. In NIPS Workshop on Practical Applications of Sparse Modeling: Open Issues and New Directions. 2010

Other contributions

• Characterization of the local minima of sparse coding

- Work in progress with R. Gribonval and F. Bach
- Non-asymptotic analysis of existence of local minima
- Extend previous work limited to under-complete and/or noiseless signals [Gribonval and Schnass, 2010; Geng et al., 2011]
- Monographs about optimization for sparse models: Related publications:

F. Bach, <u>R. Jenatton</u>, J. Mairal and G. Obozinski. Convex Optimization with Sparsity-Inducing Norms. In S. Sra, S. Nowozin, and S. J. Wright, editors, *Optimization for Machine Learning*, 2011.
F. Bach, <u>R. Jenatton</u>, J. Mairal and G. Obozinski. Optimization with Sparsity-Inducing Penalties. *To* appear in Foundations and Trends in Machine Learning, 2011.

Efficient implementations for

- Sparse structured PCA
- Active-set algorithms
- Available at www.di.ens.fr/~jenatton/

Part I Structured sparsity-inducing norms: A guided tour

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Traditional sparsity-inducing penalties

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[f(\mathbf{w})+\lambda\Omega(\mathbf{w})\right]$$

- f, usually a convex data-fitting term, e.g., square loss
- Ω , a regularization that encourages sparse solutions
- $\lambda \geq 0$, regularization parameter

Traditional sparsity-inducing penalties

$$\min_{\mathbf{w} \in \mathbb{R}^p} \left[f(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \right]$$

- f, usually a convex data-fitting term, e.g., square loss
- Ω , a regularization that encourages sparse solutions
- $\lambda \geq$ 0, regularization parameter
- Natural way to express sparsity

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_0 = \Big|\{j \in \llbracket 1; p \rrbracket; \ \mathbf{w}_j \neq 0\}\Big|$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

• Leads to combinatorial problems [Natarajan, 1995]

Traditional sparsity-inducing penalties

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[f(\mathbf{w})+\lambda\Omega(\mathbf{w})\right]$$

- f, usually a convex data-fitting term, e.g., square loss
- Ω , a regularization that encourages sparse solutions
- $\lambda \geq 0$, regularization parameter
- Natural way to express sparsity

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_0 = \left| \{j \in \llbracket 1; p \rrbracket; \ \mathbf{w}_j \neq 0 \} \right|$$

- Leads to combinatorial problems [Natarajan, 1995]
- Convex relaxation via the l₁-norm:

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_1 = \sum_{j=1}^p |\mathbf{w}_j|$$

 For least-squares regression, Lasso [Tibshirani, 1996] and basis pursuit [Chen et al., 1998]

- ℓ_0 and ℓ_1 : invariant w.r.t. permutations of ${\boldsymbol w}$
- Natural extension when variables can be grouped:

$$\Omega(\mathbf{w}) = \sum_{g \in \mathcal{G}} \left[\sum_{j \in g} \mathbf{w}_j^2 \right]^{1/2} = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2, \text{ with } \mathcal{G} \text{ a partition of } \llbracket 1; p \rrbracket$$

• Mixed ℓ_1/ℓ_2 -norm [Turlach et al., 2005; Yuan and Lin, 2006]

• Can be extended to ℓ_1/ℓ_q -norms, for $q\in(1,\infty]$

- ℓ_0 and ℓ_1 : invariant w.r.t. permutations of ${\boldsymbol w}$
- Natural extension when variables can be grouped:

$$\Omega(\mathbf{w}) = \sum_{g \in \mathcal{G}} \left[\sum_{j \in g} \mathbf{w}_j^2 \right]^{1/2} = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2, \text{ with } \mathcal{G} \text{ a partition of } \llbracket 1; p \rrbracket$$

- Mixed ℓ_1/ℓ_2 -norm [Turlach et al., 2005; Yuan and Lin, 2006]
- Can be extended to ℓ_1/ℓ_q -norms, for $q\in(1,\infty]$
- Put to zero entire groups of variables $g \in \mathcal{G}$



- ℓ_0 and $\ell_1:$ invariant w.r.t. permutations of \boldsymbol{w}
- Natural extension when variables can be grouped:

$$\Omega(\mathbf{w}) = \sum_{g \in \mathcal{G}} \left[\sum_{j \in g} \mathbf{w}_j^2 \right]^{1/2} = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2, \text{ with } \mathcal{G} \text{ a partition of } \llbracket 1; p \rrbracket$$

- Mixed ℓ_1/ℓ_2 -norm [Turlach et al., 2005; Yuan and Lin, 2006]
- Can be extended to ℓ_1/ℓ_q -norms, for $q\in(1,\infty]$
- Statistical gain when relevant prior knowledge
 - E.g., Stojnic et al. [2009]; Huang and Zhang [2010]

- ℓ_0 and $\ell_1:$ invariant w.r.t. permutations of \boldsymbol{w}
- Natural extension when variables can be grouped:

$$\Omega(\mathbf{w}) = \sum_{g \in \mathcal{G}} \left[\sum_{j \in g} \mathbf{w}_j^2 \right]^{1/2} = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2, \text{ with } \mathcal{G} \text{ a partition of } \llbracket 1; p \rrbracket$$

- Mixed ℓ_1/ℓ_2 -norm [Turlach et al., 2005; Yuan and Lin, 2006]
- Can be extended to ℓ_1/ℓ_q -norms, for $q\in(1,\infty]$
- Statistical gain when relevant prior knowledge
 - E.g., Stojnic et al. [2009]; Huang and Zhang [2010]
- Applications:
 - Encoding categorical variables [Roth and Fischer, 2008]

• Multitask learning [Obozinski et al., 2009]

Overlapping groups of variables

- Main idea: \mathcal{G} not a partition anymore
 - Ω is still a norm, $\Omega(\mathbf{w}) = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2$
- Intuitively, same behavior:
 - Some groups $g \in \mathcal{G}$ will be zeroed out
 - Sets of possible zero patterns and possible non-zero patterns?

Overlapping groups of variables

- Main idea: ${\mathcal G}$ not a partition anymore
 - Ω is still a norm, $\Omega(\mathbf{w}) = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2$
- Intuitively, same behavior:
 - Some groups $g \in \mathcal{G}$ will be zeroed out
 - Sets of possible zero patterns and possible non-zero patterns?

$$\mathcal{Z} = \Big\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \Big\} \quad \text{and} \quad \mathcal{N} = \Big\{ \bigcap_{g \in \mathcal{G}'} g^c; \ \mathcal{G}' \subseteq \mathcal{G} \Big\}$$

Overlapping groups of variables

- Main idea: ${\mathcal G}$ not a partition anymore
 - Ω is still a norm, $\Omega(\mathbf{w}) = \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2$
- Intuitively, same behavior:
 - Some groups $g \in \mathcal{G}$ will be zeroed out
 - Sets of possible zero patterns and possible non-zero patterns?

$$\mathcal{Z} = \Big\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \Big\} \quad \text{and} \quad \mathcal{N} = \Big\{ \bigcap_{g \in \mathcal{G}'} g^c; \ \mathcal{G}' \subseteq \mathcal{G} \Big\}$$



Example: Selection of contiguous patterns



The nonzero patterns are all the segments

Example: Selection of rectangles



The nonzero patterns are all the rectangles

Example: Hierarchical structure

[Zhao et al., 2009]



- Groups in ${\mathcal G}$ are all the rooted subtrees
- Selection rule:
 - If a node is selected, the same goes for all its ancestors
 - If a node is not selected, then its descendants are not selected

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

$$\mathcal{Z} = \Big\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \Big\} \quad \text{and} \quad \mathcal{N} = \Big\{ \bigcap_{g \in \mathcal{G}'} g^c; \ \mathcal{G}' \subseteq \mathcal{G} \Big\}$$

- Set of zero patterns ${\mathcal Z}$ closed under union
- Set of nonzero patterns ${\cal N}$ closed under intersection

$$\mathcal{Z} = \Big\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \Big\} \quad \text{and} \quad \mathcal{N} = \Big\{ \bigcap_{g \in \mathcal{G}'} g^c; \ \mathcal{G}' \subseteq \mathcal{G} \Big\}$$

- Set of zero patterns \mathcal{Z} closed under union
- Set of nonzero patterns ${\mathcal N}$ closed under intersection
- Relation between \mathcal{G} and \mathcal{Z} , \mathcal{N} ?

$$\mathcal{Z} = \Big\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \Big\} \quad \text{and} \quad \mathcal{N} = \Big\{ \bigcap_{g \in \mathcal{G}'} g^c; \ \mathcal{G}' \subseteq \mathcal{G} \Big\}$$

- Set of zero patterns \mathcal{Z} closed under union
- Set of nonzero patterns ${\mathcal N}$ closed under intersection
- Relation between \mathcal{G} and \mathcal{Z} , \mathcal{N} ?
- From \mathcal{G} to \mathcal{Z} :
 - Union closure of ${\mathcal G}$

$$\mathcal{Z} = \Big\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \Big\} \quad \text{and} \quad \mathcal{N} = \Big\{ \bigcap_{g \in \mathcal{G}'} g^c; \ \mathcal{G}' \subseteq \mathcal{G} \Big\}$$

- Set of zero patterns \mathcal{Z} closed under union
- Set of nonzero patterns ${\cal N}$ closed under intersection
- Relation between \mathcal{G} and \mathcal{Z} , \mathcal{N} ?
- From \mathcal{G} to \mathcal{Z} :
 - Union closure of ${\cal G}$
- From (union-closed) \mathcal{Z} to (minimal) \mathcal{G} :
 - Procedure from set theory [Doignon and Falmagne, 1998]
 - Automatic norm design

Some theoretical properties

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n I(y_i,\mathbf{w}^{\top}\mathbf{x}_i) + \lambda\Omega(\mathbf{w})\right] \quad (*)$$

- Supervised setting
 - $\{(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}\}_{i \in [\![1:n]\!]}$, (input, output) data points
 - loss function *I*, convex, \tilde{C}^2
Some theoretical properties

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n I(y_i,\mathbf{w}^{\top}\mathbf{x}_i) + \lambda\Omega(\mathbf{w})\right] \quad (*)$$

- Supervised setting
 - $\{(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}\}_{i \in \llbracket 1; n \rrbracket}$, (input, output) data points
 - loss function I, convex, C^2
- Zero patterns:
 - Solutions of (*) have zero patterns in $\mathcal{Z} = \left\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \right\}$

Some theoretical properties

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n I(y_i,\mathbf{w}^{\top}\mathbf{x}_i) + \lambda\Omega(\mathbf{w})\right] \quad (*)$$

- Supervised setting
 - $\{(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}\}_{i \in \llbracket 1; n \rrbracket}$, (input, output) data points
 - loss function I, convex, C^2
- Zero patterns:
 - Solutions of (*) have zero patterns in $\mathcal{Z} = \left\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \right\}$
- Support recovery for least-squares regression:
 - For ℓ_1 [Zhao and Yu, 2006] and group Lasso [Bach, 2008]
 - Low- and high-dimensional settings

Related approaches to structured sparsity

• Convex approaches

- Union-closed nonzero patterns [Jacob et al., 2009]
- Via convex cones [Micchelli et al., 2010]
- Submodular approaches
 - Set-functions based on supports and level-sets [Bach, 2010a,b]

Non-convex approaches

- Information-theoretic criterion [Haupt and Nowak, 2006; Huang et al., 2009]
- Union of supports [Baraniuk et al., 2010]
- Bayesian approaches
 - Hierarchical wavelet model [He and Carin, 2009]

Part II

Optimization for sparsity-inducing norms: The case of tree-structured groups

▲ロト ▲圖 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ● 臣 ● のへで

Which optimization tools?

$$\min_{\mathbf{w}\in\mathbb{R}^{p}}\left[f(\mathbf{w})+\lambda\Omega(\mathbf{w})\right]$$

- f convex, differentiable, Lipschitz-continuous gradient
- Ω convex, but nonsmooth
- Techniques blind to composite structure:
 - Subgradient descent
 - Interior-point for QP, SOCP

Which optimization tools?

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[f(\mathbf{w})+\lambda\Omega(\mathbf{w})\right]$$

- f convex, differentiable, Lipschitz-continuous gradient
- Ω convex, but nonsmooth
- Techniques blind to composite structure:
 - Subgradient descent
 - Interior-point for QP, SOCP
- Techniques developed in the thesis:
 - Active-set strategies
 - Reweighted- ℓ_2 schemes
 - Proximal methods

Proximal methods

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- Forward-backward splitting [Combettes and Pesquet, 2010]
- At iteration k, update rule

$$\mathbf{w}^{(k+1)} \leftarrow \arg\min_{\mathbf{w}\in\mathbb{R}^{p}} \left[f(\mathbf{w}^{(k)}) + \nabla f(\mathbf{w}^{(k)})^{\top} (\mathbf{w} - \mathbf{w}^{(k)}) + \frac{L}{2} \|\mathbf{w}^{(k)} - \mathbf{w}\|_{2}^{2} + \lambda \Omega(\mathbf{w}) \right]$$

$$\leftarrow \arg\min_{\mathbf{w}\in\mathbb{R}^{p}} \left[\frac{1}{2} \left\| \mathbf{w} - (\mathbf{w}^{(k)} - \frac{1}{L} \nabla f(\mathbf{w}^{(k)})) \right\|_{2}^{2} + \lambda \Omega(\mathbf{w}) \right]$$

- Linearization of f at **w**^(k)
- *Proximity* term, remain around $\mathbf{w}^{(k)}$

Proximal methods

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

- Forward-backward splitting [Combettes and Pesquet, 2010]
- At iteration k, update rule

$$\mathbf{w}^{(k+1)} \leftarrow \arg\min_{\mathbf{w}\in\mathbb{R}^{p}} \left[f(\mathbf{w}^{(k)}) + \nabla f(\mathbf{w}^{(k)})^{\top} (\mathbf{w} - \mathbf{w}^{(k)}) + \frac{L}{2} \|\mathbf{w}^{(k)} - \mathbf{w}\|_{2}^{2} + \lambda \Omega(\mathbf{w}) \right]$$

$$\leftarrow \arg\min_{\mathbf{w}\in\mathbb{R}^{p}} \left[\frac{1}{2} \left\| \mathbf{w} - (\mathbf{w}^{(k)} - \frac{1}{L} \nabla f(\mathbf{w}^{(k)})) \right\|_{2}^{2} + \lambda \Omega(\mathbf{w}) \right]$$

- Linearization of f at **w**^(k)
- *Proximity* term, remain around **w**^(k)
- If $\lambda = 0$, gradient update rule: $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} \frac{1}{L} \nabla f(\mathbf{w}^{(k)})$

Proximal methods

- Forward-backward splitting [Combettes and Pesquet, 2010]
- At iteration k, update rule

$$\mathbf{w}^{(k+1)} \leftarrow \arg\min_{\mathbf{w}\in\mathbb{R}^{p}} \left[f(\mathbf{w}^{(k)}) + \nabla f(\mathbf{w}^{(k)})^{\top} (\mathbf{w} - \mathbf{w}^{(k)}) + \frac{L}{2} \|\mathbf{w}^{(k)} - \mathbf{w}\|_{2}^{2} + \lambda \Omega(\mathbf{w}) \right]$$

$$\leftarrow \arg\min_{\mathbf{w}\in\mathbb{R}^{p}} \left[\frac{1}{2} \left\| \mathbf{w} - (\mathbf{w}^{(k)} - \frac{1}{L} \nabla f(\mathbf{w}^{(k)})) \right\|_{2}^{2} + \lambda \Omega(\mathbf{w}) \right]$$

- Linearization of f at **w**^(k)
- *Proximity* term, remain around $\mathbf{w}^{(k)}$
- If $\lambda = 0$, gradient update rule: $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} \frac{1}{L} \nabla f(\mathbf{w}^{(k)})$
- There are also accelerated versions [Beck and Teboulle, 2009; Nesterov, 2007]
 - Linear combination of past iterates
 - Optimal convergence rates $O(\frac{1}{k^2})$ [Nesterov, 2004]

• Idea: Ω only involved through its proximal operator:

$$\mathsf{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \argmin_{\mathbf{w} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \Omega(\mathbf{w}) \right]$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

• Has to be solved efficiently and exactly [or with arbitrary precision, Schmidt et al., 2011]

• Idea: Ω only involved through its proximal operator:

$$\mathsf{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \argmin_{\mathbf{w} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \Omega(\mathbf{w}) \right]$$

- Has to be solved efficiently and exactly [or with arbitrary precision, Schmidt et al., 2011]
- For ℓ_1 -norm, soft-thresholding [Donoho and Johnstone, 1995]:

$$\left[\mathsf{Prox}_{\lambda\parallel,\parallel_1}(\mathbf{u})\right]_j = \operatorname{sign}(\mathbf{u}_j) \max\{0, |\mathbf{u}_j| - \lambda\}$$

• Idea: Ω only involved through its proximal operator:

$$\mathsf{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \argmin_{\mathbf{w} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \Omega(\mathbf{w}) \right]$$

- Has to be solved efficiently and exactly [or with arbitrary precision, Schmidt et al., 2011]
- For ℓ_1 -norm, soft-thresholding [Donoho and Johnstone, 1995]:

$$\left[\mathsf{Prox}_{\lambda\parallel,\parallel_1}(\mathbf{u})\right]_j = \operatorname{sign}(\mathbf{u}_j) \max\{0, |\mathbf{u}_j| - \lambda\}$$

• Also simple expressions for ℓ_2 - and ℓ_∞ -norms

• Idea: Ω only involved through its proximal operator:

$$\mathsf{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \argmin_{\mathbf{w} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \Omega(\mathbf{w}) \right]$$

- Has to be solved efficiently and exactly [or with arbitrary precision, Schmidt et al., 2011]
- For ℓ_1 -norm, soft-thresholding [Donoho and Johnstone, 1995]:

$$\left[\mathsf{Prox}_{\lambda\parallel,\parallel_1}(\mathbf{u})\right]_j = \operatorname{sign}(\mathbf{u}_j) \max\{0, |\mathbf{u}_j| - \lambda\}$$

- Also simple expressions for ℓ_2 and ℓ_∞ -norms
- When Ω has overlapping groups, no closed-form available...

Tree-structured set of groups



• Tree-structured set of groups \mathcal{G} :

For any
$$g, h \in \mathcal{G}$$
, $\left[g \cap h \neq \varnothing\right] \Rightarrow \left[g \subseteq h \text{ or } h \subseteq g\right]$

Proximal operator for tree-structured $\mathcal G$

$$\mathsf{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \argmin_{\mathbf{w} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_q \right] \quad \text{for } q \in \{2, \infty\}$$

Proximal operator for tree-structured $\mathcal G$

$$\operatorname{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \operatorname*{arg\,min}_{\mathbf{w} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_q \right] \quad \text{for } q \in \{2, \infty\}$$

Consider

$$\operatorname{Prox}_{g}: \ \mathbf{u} \mapsto \operatorname*{arg\,min}_{\mathbf{w} \in \mathbb{R}^{p}} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_{2}^{2} + \lambda \|\mathbf{w}_{g}\|_{q} \right] \quad \text{for } q \in \{2, \infty\}$$

Proximal operator for tree-structured \mathcal{G}

$$\operatorname{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \operatorname*{arg\,min}_{\mathbf{w} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_2^2 + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_q \right] \quad \text{for } q \in \{2, \infty\}$$

Consider

$$\operatorname{Prox}_{g}: \mathbf{u} \mapsto \operatorname*{arg\,min}_{\mathbf{w} \in \mathbb{R}^{p}} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_{2}^{2} + \lambda \|\mathbf{w}_{g}\|_{q} \right] \quad \text{for } q \in \{2, \infty\}$$

- Assume $\mathcal{G} = \{g_1, \dots, g_{|\mathcal{G}|}\}$ ordered according to

$$g \preceq h \Rightarrow \Big[g \subseteq h \text{ or } g \cap h = \varnothing\Big]$$
 ("leaves up to the root")

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Proximal operator for tree-structured $\mathcal G$

$$\operatorname{Prox}_{\lambda\Omega}: \ \mathbf{u} \mapsto \underset{\mathbf{w} \in \mathbb{R}^{p}}{\operatorname{arg\,min}} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_{2}^{2} + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{w}_{g}\|_{q} \right] \quad \text{for } q \in \{2, \infty\}$$

Consider

$$\operatorname{Prox}_{g}: \mathbf{u} \mapsto \operatorname*{arg\,min}_{\mathbf{w} \in \mathbb{R}^{p}} \left[\frac{1}{2} \|\mathbf{w} - \mathbf{u}\|_{2}^{2} + \lambda \|\mathbf{w}_{g}\|_{q} \right] \quad \text{for } q \in \{2, \infty\}$$

• Assume
$$\mathcal{G} = \{g_1, \dots, g_{|\mathcal{G}|}\}$$
 ordered according to

$$g \preceq h \Rightarrow \left\lfloor g \subseteq h \text{ or } g \cap h = \varnothing
ight
brace$$
 ("leaves up to the root")

Proposition (composition of prox)

$$Prox_{\lambda\Omega} = Prox_{g_{|\mathcal{G}|}} \circ \cdots \circ Prox_{g_1}$$

With careful implementation, complexity linear O(p) when q = 2.

Extensions

◆□▶ ◆□▶ ▲□▶ ▲□▶ ▲□ ◆ ●

- General overlapping groups for ℓ_1/ℓ_∞ -norms
 - Mairal, Jenatton, Obozinski, and Bach [2011]
 - Connections with network flow problems
 - Conic duality + divide-and-conquer strategy

Part III Some applications to dictionary learning

Dictionary learning in a nutshell

- **Data**: *n* signals, $\mathbf{X} = [\mathbf{x}^1, \dots, \mathbf{x}^n] \in \mathbb{R}^{m \times n}$
- Dictionary: p atoms, $\mathbf{D} = [\mathbf{d}^1, \dots, \mathbf{d}^p] \in \mathbb{R}^{m \times p}$
- **Decomposition**: $\mathbf{A} = [\boldsymbol{\alpha}^1, \dots, \boldsymbol{\alpha}^n] \in \mathbb{R}^{p \times n}$
- Goal:

$\mathbf{X} \approx \mathbf{D}\mathbf{A}$

- Applied to various settings and classes of signals
 - Neuroscience [Olshausen and Field, 1997]
 - Image processing/Computer vision [Elad and Aharon, 2006; Peyré, 2009; Mairal, 2010]
 - Audio processing [Sprechmann et al., 2010; Lefèvre et al., 2011]

- Topic modeling [Jenatton et al., 2011]
- . . .

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

• Mainly, two settings

- Mainly, two settings
- Regularization over **D**: $\Omega(\mathbf{d}^{j})$, for $j \in [\![1; p]\!]$
 - Sparse structured dictionary elements



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

- Regularization over A: $\Omega(\alpha^i)$, for $i \in \llbracket 1; n \rrbracket$
 - Each signal uses few atoms, sparse decomposition
 - Induced effect on D
 - Atoms organize and adapt to match the structure of $\boldsymbol{\Omega}$



▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- Regularization over A: $\Omega(\alpha^i)$, for $i \in \llbracket 1; n \rrbracket$
 - Each signal uses few atoms, sparse decomposition
 - Induced effect on D
 - Atoms organize and adapt to match the structure of $\boldsymbol{\Omega}$



Example, the hierarchical case:

$$\min_{\mathbf{A}\in\mathbb{R}^{p\times n},\mathbf{D}\in\mathcal{D}}\frac{1}{n}\sum_{i=1}^{n}\left[\frac{1}{2}\|\mathbf{x}^{i}-\mathbf{D}\boldsymbol{\alpha}^{i}\|_{2}^{2}+\lambda\Omega_{\text{hierarchical}}(\boldsymbol{\alpha}^{i})\right]$$



- Graph structure imposed by Ω
- Each patch at a node represents an atom

Part IV Application to neuroimaging

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Structured sparsity for neuroimaging

- Supervised problem
- $\mathbf{x} \in \mathbb{R}^{p}$ fMRI signal, with $p \approx 70\,000$ voxels
- y sizes of objects
- Goal: Predict y from brain activation-map x, across subjects

Size *y* of the object?

fMRI signal ${\bf x}$



Predict



< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

- 3

Structured sparsity for neuroimaging

- Supervised problem
- $\mathbf{x} \in \mathbb{R}^{p}$ fMRI signal, with $p \approx 70\,000$ voxels
- y sizes of objects
- Goal: Predict y from brain activation-map x, across subjects
- Some properties:
 - High-dimensional problem, only $n \approx 100$
 - A few voxels are useful to predict y (sparsity)
 - Relevant voxels appear as spatially-localized patterns

- Across subjects, voxel-based predictions not reliable:
 - Misalignement issue

Proposed approach

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

э

500

• Instead of single voxels, consider local averages: parcels



Proposed approach

• Instead of single voxels, consider local averages: parcels



• Build a multi-scale representation of the data:

Proposed approach

Instead of single voxels, consider local averages: parcels



- Build a multi-scale representation of the data:
 - Ward's hierarchical clustering + spatial constraints
 - Leaves of the tree: single voxels
 - Root of the tree: all voxels

The new representation of the signals



Original signal

$$\mathbf{x}_{old} = [\mathbf{x}_{voxel1}, \mathbf{x}_{voxel2}, \mathbf{x}_{voxel3}] \in \mathbb{R}^3$$

・ロト ・四ト ・ヨト ・ヨト 三日

The new representation of the signals



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

Reminiscent of spatial pyramid from Lazebnik et al. [2006]

How to exploit the new representation?



- $\mathbf{x}_{new} = [\mathbf{x}_{old}, \mathbf{x}_{parcel1}, \mathbf{x}_{parcel2}] \in \mathbb{R}^5$
- Assume linear models: $\mathbf{w}^{\top}\mathbf{x}_{new}$
- Idea: The closer to the leaves, the less reliable across subjects

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

• First select **w**_j near the root

How to exploit the new representation?



- $\mathbf{x}_{new} = [\mathbf{x}_{old}, \mathbf{x}_{parcel1}, \mathbf{x}_{parcel2}] \in \mathbb{R}^5$
- Assume linear models: $\mathbf{w}^{\top}\mathbf{x}_{new}$
- Idea: The closer to the leaves, the less reliable across subjects

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

- First select **w**_j near the root
- ℓ_1 -norm would not distinguish between parcels/voxels
How to exploit the new representation?



- $\mathbf{x}_{new} = [\mathbf{x}_{old}, \mathbf{x}_{parcel1}, \mathbf{x}_{parcel2}] \in \mathbb{R}^5$
- Assume linear models: $\mathbf{w}^{\top} \mathbf{x}_{new}$
- Idea: The closer to the leaves, the less reliable across subjects

- First select **w**_j near the root
- ℓ_1 -norm would not distinguish between parcels/voxels
- We use a hierarchical sparsity-inducing norm

How to exploit the new representation?



- $\mathbf{x}_{new} = [\mathbf{x}_{old}, \mathbf{x}_{parcel1}, \mathbf{x}_{parcel2}] \in \mathbb{R}^5$
- Assume linear models: $\mathbf{w}^{\top}\mathbf{x}_{new}$
- Idea: The closer to the leaves, the less reliable across subjects

- First select **w**_j near the root
- ℓ_1 -norm would not distinguish between parcels/voxels
- We use a hierarchical sparsity-inducing norm
 - Cannot select a voxel without selecting ancestral parcels

How to exploit the new representation?



- $\mathbf{x}_{new} = [\mathbf{x}_{old}, \mathbf{x}_{parcel1}, \mathbf{x}_{parcel2}] \in \mathbb{R}^5$
- Assume linear models: $\mathbf{w}^{\top}\mathbf{x}_{new}$
- Idea: The closer to the leaves, the less reliable across subjects
- First select w_j near the root
- ℓ_1 -norm would not distinguish between parcels/voxels
- We use a hierarchical sparsity-inducing norm
 - Cannot select a voxel without selecting ancestral parcels
 - If a parcel is zeroed out, the same goes for all its descendants

Regression results

Loss function:	Square	
	Error (mean,std)	P-value w.r.t. Tree ℓ_2
Regularization:		
ℓ_2 (Ridge)	(13.8, 7.6)	0.096
ℓ_1	(20.2, 10.9)	0.013*
$\ell_1 + \ell_2$ (Elastic net)	(14.4, 8.8)	0.065
Reweighted ℓ_1	(18.8, 14.6)	0.052
ℓ_1 (augmented space)	(14.2, 7.9)	0.096
ℓ_1 (tree weights)	(13.9, 7.9)	0.032*
Tree ℓ_2	(11.8, 6.7)	-
Tree ℓ_∞	(12.8, 6.7)	0.137
Greedy [Michel et al., 2010]	(12.0, 5.5)	0.5

Classification results

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Loss function:	Multinomial logistic	
	Error (mean,std)	P-value w.r.t. Tree ℓ_2
Regularization:		
ℓ_2 (Ridge)	(24.2, 9.2)	0.035*
ℓ_1	(25.8, 12.0)	0.004*
ℓ_1/ℓ_2 (Multi-task)	(26.7, 7.6)	0.007*
ℓ_1/ℓ_∞ (Multi-task)	(26.7, 11.6)	0.002*
Tree ℓ_2	(16.7, 10.4)	-
Tree ℓ_∞	(22.5, 13.0)	0.156
Greedy [Michel et al., 2010]	(21.6, 14.5)	0.001*

Neuroimaging conclusions

- Best prediction accuracy
- Multi-scale representation
 - Robust to misalignment
 - Adapted to the spatial encoding of brain activations
- Explored in a convex way
 - No initialization issue
 - Minimum guarantee
- Improve supervised learning from brain imaging data:
 - fMRI decoding like shown here
 - · Ability to classify patients vs. controls for diagnosis purposes

Conclusions and take-home messages

- Structured sparsity via convex tools
- Structure ⇔ Constraints on possible sparsity patterns
- Plain and structured sparsity for the same computational cost!
- Different scenarios where structured sparsity applies:
 - Fixed dictionary with *preprocessing* to match Ω

 -e.g., neuroimaging with hierarchical clustering
 - Dictionary learning with learned atoms matching Ω -e.g., tree of image patches

- Connections between:
 - Sparse structured dictionary learning
 - Non-parametric Bayesian techniques

Perspectives and future work

- Deeper statistical analysis
 - Quantify statistical gain w.r.t. ℓ_1 -norm
 - E.g., Negahban et al. [2009]; Maurer and Pontil [2011]
- Understand the role of regularizers in dictionary learning
 - Risk analysis [Maurer and Pontil, 2010]
 - Identifiability questions [Gribonval and Schnass, 2010]
- Monitor better the optimization in dictionary learning
- Further validate structured sparsity in NLP
 - Structures are plentiful [Martins et al., 2011]

Thank you all for your attention

<□ > < @ > < E > < E > E のQ @

References I

- F. Bach. Consistency of the group lasso and multiple kernel learning. *Journal of Machine Learning Research*, 9:1179–1225, 2008.
- F. Bach. Structured sparsity-inducing norms through submodular functions. In *Advances in Neural Information Processing Systems*, 2010a.
- F. Bach. Shaping level sets with submodular functions. Technical report, Preprint arXiv:1012.1501, 2010b.
- R. G. Baraniuk, V. Cevher, M. F. Duarte, and C. Hegde. Model-based compressive sensing. *IEEE Transactions on Information Theory*, 56: 1982–2001, 2010.
- A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM Journal on Imaging Sciences, 2(1):183–202, 2009.
- D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- D. Blei, T. L. Griffiths, and M. I. Jordan. The nested chinese restaurant process and bayesian nonparametric inference of topic hierarchies. *Journal of the ACM*, 57(2):1–30, 2010.

References II

- Y. Boykov, O. Veksler, and R. Zabih. Fast approximate energy minimization via graph cuts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(11):1222–1239, 2001.
- E.J. Candes and T. Tao. Decoding by linear programming. *IEEE Transactions* on Information Theory, 51(12):4203–4215, 2005.
- S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. SIAM Journal on Scientific Computing, 20(1):33–61, 1998.
- D. B. Chklovskii and A. A. Koulakov. Maps in the brain: What can we learn from them? *Annual Reviews in Neuroscience*, 27(1):369, 2004.
- P. L. Combettes and J.-C. Pesquet. Proximal splitting methods in signal processing. In *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*. Springer, 2010.
- J. P. Doignon and J. C. Falmagne. Knowledge Spaces. Springer-Verlag, 1998.
- D. L. Donoho and I. M. Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association*, 90(432), 1995.

References III

- M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image Processing*, 54(12):3736–3745, December 2006.
- Q. Geng, H. Wang, and J. Wright. On the Local Correctness of L1 Minimization for Dictionary Learning. Technical report, Preprint arXiv:1101.5672, 2011.
- A. Gramfort, D. Strohmeier, J. Haueisen, M. Hamalainen, and M. Kowalski. Functional brain imaging with M/EEG using structured sparsity in time-frequency dictionaries. In *Information Processing in Medical Imaging*, pages 600–611. Springer, 2011.
- R. Gribonval and K. Schnass. Dictionary identification—sparse matrix-factorization via ℓ_1 -minimization. *IEEE Transactions on Information Theory*, 56(7):3523–3539, 2010.
- J. Haupt and R. Nowak. Signal reconstruction from noisy random projections. *IEEE Transactions on Information Theory*, 52(9):4036–4048, 2006.
- L. He and L. Carin. Exploiting structure in wavelet-based Bayesian compressive sensing. *IEEE Transactions on Signal Processing*, 57:3488–3497, 2009.
- J. Huang and T. Zhang. The benefit of group sparsity. *Annals of Statistics*, 38 (4):1978–2004, 2010.

References IV

- J. Huang, T. Zhang, and D. Metaxas. Learning with structured sparsity. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2009.
- L. Jacob, G. Obozinski, and J.-P. Vert. Group Lasso with overlaps and graph Lasso. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2009.
- R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for hierarchical sparse coding. *Journal of Machine Learning Research*, 12: 2297–2334, 2011.
- S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *Proceedings* of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2006.
- A. Lefèvre, F. Bach, and C. Févotte. Itakura-saito nonnegative matrix factorization with group sparsity. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2011.
- J. Mairal. Sparse coding for machine learning, image processing and computer vision. PhD thesis, École normale supérieure de Cachan ENS Cachan, 2010. Available at http://tel.archives-ouvertes.fr/tel-00595312/fr/.

References V

- J. Mairal, R. Jenatton, G. Obozinski, and F. Bach. Convex and network flow optimization for structured sparsity. *Journal of Machine Learning Research*, 12:2681–2720, 2011.
- A. F. T. Martins, N. A. Smith, P. M. Q. Aguiar, and M. A. T. Figueiredo. Structured sparsity in structured prediction. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2011.
- A. Maurer and M. Pontil. *k*-dimensional coding schemes in hilbert spaces. *IEEE Transactions on Information Theory*, 56(11):5839–5846, 2010.
- A. Maurer and M. Pontil. Structured sparsity and generalization. Technical report, Preprint arXiv:1108.3476, 2011.
- N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the lasso. *Annals of Statistics*, 34(3):1436–1462, 2006.
- C. A. Micchelli, J. M. Morales, and M. Pontil. A family of penalty functions for structured sparsity. In *Advances in Neural Information Processing Systems*, 2010.
- V. Michel, E. Eger, C. Keribin, J.-B. Poline, and B. Thirion. A supervised clustering approach for extracting predictive information from brain activation images. *MMBIA'10*, 2010.

References VI

- B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM journal* on computing, 24:227, 1995.
- S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. In *Advances in Neural Information Processing Systems*, 2009.
- Y. Nesterov. *Introductory lectures on convex optimization: a basic course.* Kluwer Academic Publishers, 2004.
- Y. Nesterov. Gradient methods for minimizing composite objective function. Technical report, Center for Operations Research and Econometrics (CORE), Catholic University of Louvain, 2007.
- G. Obozinski, B. Taskar, and M. I. Jordan. Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, pages 1–22, 2009.
- B. A. Olshausen and D. J. Field. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, 37:3311–3325, 1997.
- G. Peyré. Sparse modeling of textures. *Journal of Mathematical Imaging and Vision*, 34(1):17–31, 2009.

References VII

- F. Rapaport, E. Barillot, and J.-P. Vert. Classification of arrayCGH data using fused SVM. *Bioinformatics*, 24(13):i375–i382, Jul 2008.
- T. Rebafka, C. Lévy-Leduc, and M. Charbit. Omp-type algorithm with structured sparsity patterns for multipath radar signals. Technical report, Preprint arXiv:1103.5158, 2011.
- V. Roth and B. Fischer. The group-Lasso for generalized linear models: uniqueness of solutions and efficient algorithms. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.
- M. Schmidt and K. Murphy. Convex structure learning in log-linear models: Beyond pairwise potentials. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2010.
- M. Schmidt, N. Le Roux, and F. Bach. Convergence rates of inexact proximal-gradient methods for convex optimization. In *Advances in Neural Information Processing Systems*, 2011.
- J. M. Shapiro. Embedded image coding using zerotrees of wavelet coefficients. *IEEE Transactions on Signal Processing*, 41(12):3445–3462, 1993.
- P. Sprechmann, I. Ramirez, P. Cancela, and G. Sapiro. Collaborative sources identification in mixed signals via hierarchical sparse modeling. Technical report, Preprint arXiv:1010.4893, 2010.

References VIII

- M. Stojnic, F. Parvaresh, and B. Hassibi. On the reconstruction of block-sparse signals with an optimal number of measurements. *IEEE Transactions on Signal Processing*, 57(8):3075–3085, 2009.
- A. Tibau Puig, A. Wiesel, A. Zaas, C. Woods, G. Ginsburg, G. Fleury, and A. Hero. Order-preserving factor analysis-application to longitudinal gene expression. *IEEE Transactions on Signal Processing*, 99(99):1–1, 2011.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B*, pages 267–288, 1996.
- B. A. Turlach, W. N. Venables, and S. J. Wright. Simultaneous variable selection. *Technometrics*, 47(3):349–363, 2005.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society. Series B*, 68(1):49–67, 2006.
- P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2563, 2006.
- P. Zhao, G. Rocha, and B. Yu. The composite absolute penalties family for grouped and hierarchical variable selection. *Annals of Statistics*, 37(6A): 3468–3497, 2009.

Part V Appendix

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n(y_i-\mathbf{w}^{\top}\mathbf{x}_i)^2+\lambda\Omega(\mathbf{w})\right]$$

- Fixed-design model: $y = \mathbf{x}^{\top} \bar{\mathbf{w}} + \varepsilon$
- $\lim_{n\to\infty} \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i \mathbf{x}_i^{\top} = \bar{\mathbf{Q}} \succ 0$
- Goal: Recover $\overline{J} = \{j \in \llbracket 1; p \rrbracket; \ \overline{\mathbf{w}}_j \neq 0\}$ by its estimate \widehat{J} :

$$\mathbb{P}(\hat{\mathrm{J}}=ar{\mathrm{J}}) \underset{n o +\infty}{\longrightarrow} 1 \quad (*)$$

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n(y_i-\mathbf{w}^{\top}\mathbf{x}_i)^2+\lambda\Omega(\mathbf{w})\right]$$

- Fixed-design model: $y = \mathbf{x}^{\top} \bar{\mathbf{w}} + \varepsilon$
- $\lim_{n\to\infty} \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i \mathbf{x}_i^{\top} = \bar{\mathbf{Q}} \succ 0$
- Goal: Recover $\overline{J} = \{j \in \llbracket 1; p \rrbracket; \ \overline{\mathbf{w}}_j \neq 0\}$ by its estimate \widehat{J} :

$$\mathbb{P}(\hat{\mathrm{J}}=ar{\mathrm{J}}) \underset{n o +\infty}{ o} \mathbb{1} \quad (*)$$

• Key quantity:

Correlation bewteen "good" and "bad" variables $\theta_{\text{struct}} = \underbrace{\Omega^*}_{\text{dual norm of }\Omega} \left(\overbrace{\bar{\bar{\mathbf{Q}}}_{\bar{\mathbf{J}}c\bar{\mathbf{J}}}}^{\infty} [\bar{\mathbf{Q}}_{\bar{\mathbf{J}}\bar{\mathbf{J}}}]^{-1} \mathbf{r}_{\bar{\mathbf{J}}} \right), \text{ with } \mathbf{r}_j = \bar{\mathbf{w}}_j \sum_{\substack{g \in \mathcal{G}, g \ni j \\ g \cap \bar{\mathbf{J}} \neq \emptyset}} \|\bar{\mathbf{w}}_g\|_2^{-1}$

▲□▶ ▲御▶ ▲臣▶ ▲臣▶ 三臣 - わへで

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n(y_i-\mathbf{w}^{\top}\mathbf{x}_i)^2+\lambda\Omega(\mathbf{w})\right]$$

- Fixed-design model: $y = \mathbf{x}^{\top} \bar{\mathbf{w}} + \varepsilon$
- $\lim_{n\to\infty} \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i \mathbf{x}_i^{\top} = \bar{\mathbf{Q}} \succ 0$
- Goal: Recover $\overline{J} = \{j \in \llbracket 1; p \rrbracket; \ \overline{\mathbf{w}}_j \neq 0\}$ by its estimate \widehat{J} :

$$\mathbb{P}(\hat{\mathrm{J}}=ar{\mathrm{J}}) \underset{n o +\infty}{\longrightarrow} 1 \quad (*)$$

Key quantity:

Correlation bewteen "good" and "bad" variables

$$\theta_{\text{struct}} = \underbrace{\Omega^*}_{\text{dual norm of }\Omega} \left(\overbrace{\bar{\mathbf{Q}}_{\bar{\mathbf{j}}c\bar{\mathbf{j}}}}^{\infty} [\bar{\mathbf{Q}}_{\bar{\mathbf{j}}\bar{\mathbf{j}}}]^{-1} \mathbf{r}_{\bar{\mathbf{j}}} \right), \text{ with } \mathbf{r}_j = \bar{\mathbf{w}}_j \sum_{\substack{g \in \mathcal{G}, g \ni j\\g \cap \bar{\mathbf{j}} \neq \emptyset}} \|\bar{\mathbf{w}}_g\|_2^{-1}$$

For the ℓ_1 -norm, $\theta_{\ell_1} = \left\| \bar{\mathbf{Q}}_{\bar{\mathbf{J}}^c \bar{\mathbf{J}}} [\bar{\mathbf{Q}}_{\bar{\mathbf{J}}\bar{\mathbf{J}}}]^{-1} \operatorname{sign}(\bar{\mathbf{w}})_{\bar{\mathbf{J}}} \right\|_{\infty}$ [Zhao and Yu, 2006]

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n(y_i-\mathbf{w}^{\top}\mathbf{x}_i)^2+\lambda\Omega(\mathbf{w})\right]$$

- Fixed-design model: $y = \mathbf{x}^{\top} \bar{\mathbf{w}} + \varepsilon$
- $\lim_{n\to\infty} \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i \mathbf{x}_i^{\top} = \bar{\mathbf{Q}} \succ 0$
- Goal: Recover $\overline{J} = \{j \in \llbracket 1; p \rrbracket; \ \overline{\mathbf{w}}_j \neq 0\}$ by its estimate \widehat{J} :

$$\mathbb{P}(\hat{\mathrm{J}}=ar{\mathrm{J}}) \underset{n o +\infty}{\longrightarrow} 1 \quad (*)$$

• Key quantity:

Correlation bewteen "good" and "bad" variables

$$\theta_{\text{struct}} = \underbrace{\Omega^*}_{\text{dual norm of }\Omega} \left(\overbrace{\bar{\mathbf{Q}}_{\bar{\mathbf{j}}c\bar{\mathbf{j}}}}^{\infty} [\bar{\mathbf{Q}}_{\bar{\mathbf{j}}\bar{\mathbf{j}}}]^{-1} \mathbf{r}_{\bar{\mathbf{j}}} \right), \text{ with } \mathbf{r}_j = \bar{\mathbf{w}}_j \sum_{\substack{g \in \mathcal{G}, g \ni j \\ g \cap \bar{\mathbf{j}} \neq \emptyset}} \|\bar{\mathbf{w}}_g\|_2^{-1}$$

Proposition

If $\lambda \to 0, \lambda \sqrt{n} \to +\infty$, then if $\theta_{struct} < 1$, then (*). Conversely, if (*), then $\theta_{struct} \le 1$. Formal characterization of zero patterns

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n l(y_i,\mathbf{w}^{\top}\mathbf{x}_i) + \lambda\Omega(\mathbf{w})\right] \quad (*)$$

• $\{(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}\}_{i \in [\![1;n]\!]}$, (input, output) data points • $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} \in \mathbb{R}^{n \times p}$ fixed, $\{y_i\}_{i \in [\![1;n]\!]}$ random (with a density)

Formal characterization of zero patterns

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n I(y_i,\mathbf{w}^\top\mathbf{x}_i) + \lambda\Omega(\mathbf{w})\right] \quad (*)$$

• $\{(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}\}_{i \in [\![1;n]\!]}$, (input, output) data points • $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} \in \mathbb{R}^{n \times p}$ fixed, $\{y_i\}_{i \in [\![1;n]\!]}$ random (with a density)

• Loss function, $I: (y, y') \mapsto I(y, y')$ convex, C^2 , with

•
$$\frac{\partial^2 I(y,y')}{\partial y \partial y'} \neq 0$$
 and $\frac{\partial^2 I(y,y')}{\partial^2 y'} > 0$

Formal characterization of zero patterns

$$\min_{\mathbf{w}\in\mathbb{R}^p}\left[\frac{1}{n}\sum_{i=1}^n I(y_i,\mathbf{w}^{\top}\mathbf{x}_i) + \lambda\Omega(\mathbf{w})\right] \quad (*)$$

- $\{(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}\}_{i \in [\![1;n]\!]}$, (input, output) data points • $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} \in \mathbb{R}^{n \times p}$ fixed, $\{y_i\}_{i \in [\![1;n]\!]}$ random (with a density)
- Loss function, $I: (y, y') \mapsto I(y, y')$ convex, C^2 , with

•
$$\frac{\partial^2 l(y,y')}{\partial y \partial y'} \neq 0$$
 and $\frac{\partial^2 l(y,y')}{\partial^2 y'} > 0$

Proposition

Let k be the column rank of **X**. Any solution of (*) with at most k - 1 nonzero entries has almost surely its zero pattern in

$$\mathcal{Z} = \Big\{ \bigcup_{g \in \mathcal{G}'} g; \ \mathcal{G}' \subseteq \mathcal{G} \Big\}.$$

(日) (同) (三) (三) (三) (○) (○)

Benchmark: Denoising of image patches



Figure: Sparsity level increasing from left to right.

- Small scale: Dictionary in $\mathbb{R}^{256 \times 151}$
- Solved thousands of times in inner loop for dictionary learning

<ロ> (四) (四) (三) (三) (三) (三)

Hierarchical topic models

- NIPS proceedings (1714 documents, 8274 words).
- Each document is modeled through word frequencies.
- Alternative to probabilistic topic models [Blei et al., 2010].

