# Structured sparse methods for matrix factorization

### Francis Bach

Sierra team, INRIA - Ecole Normale Supérieure





# March 2011 Joint work with **R. Jenatton, J. Mairal, G. Obozinski**

# Structured sparse methods for matrix factorization Outline

- Learning problems on matrices
- Sparse methods for matrices
  - Sparse principal component analysis
  - Dictionary learning
- Structured sparse PCA
  - Sparsity-inducing norms and overlapping groups
  - Structure on dictionary elements
  - Structure on decomposition coefficients

### **Learning on matrices - Collaborative filtering**

- Given  $n_{\mathcal{X}}$  "movies"  $\mathbf{x} \in \mathcal{X}$  and  $n_{\mathcal{Y}}$  "customers"  $\mathbf{y} \in \mathcal{Y}$ ,
- Predict the "rating"  $z(\mathbf{x},\mathbf{y})\in\mathcal{Z}$  of customer  $\mathbf{y}$  for movie  $\mathbf{x}$
- Training data: large  $n_X \times n_Y$  incomplete matrix **Z** that describes the known ratings of some customers for some movies
- Goal: complete the matrix.



### Learning on matrices - Image denoising

- Simultaneously denoise all patches of a given image
- Example from Mairal, Bach, Ponce, Sapiro, and Zisserman (2009c)



#### **Learning on matrices - Source separation**

• Single microphone (Benaroya et al., 2006; Févotte et al., 2009)



### Learning on matrices - Multi-task learning

- k linear prediction tasks on same covariates  $\mathbf{x} \in \mathbb{R}^p$ 
  - k weight vectors  $\mathbf{w}_j \in \mathbb{R}^p$
  - Joint matrix of predictors  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_k) \in \mathbb{R}^{p imes k}$
- Classical applications
  - Transfer learning
  - Multi-category classification (one task per class) (Amit et al., 2007)

#### • Share parameters between tasks

 Joint variable or feature selection (Obozinski et al., 2009; Pontil et al., 2007)

#### **Learning on matrices - Dimension reduction**

- Given data matrix  $\mathbf{X} = (\mathbf{x}_1^{\top}, \dots, \mathbf{x}_n^{\top}) \in \mathbb{R}^{n \times p}$ 
  - Principal component analysis:  $|\mathbf{x}_i pprox \mathbf{D} m{lpha}_i|$



# **Sparsity in machine learning**

- Assumption:  $\mathbf{y} = \mathbf{w}^{\top}\mathbf{x} + \boldsymbol{\varepsilon}$ , with  $\mathbf{w} \in \mathbb{R}^p$  sparse
  - Proxy for interpretability
  - Allow high-dimensional inference: log

$$\log p = O(n)$$

• **Sparsity and convexity** ( $\ell_1$ -norm regularization):





# Two types of sparsity for matrices $\mathbf{M} \in \mathbb{R}^{n \times p}$ I - Directly on the elements of $\mathbf{M}$

• Many zero elements:  $\mathbf{M}_{ij} = 0$ 



• Many zero rows (or columns):  $(\mathbf{M}_{i1}, \ldots, \mathbf{M}_{ip}) = 0$ 



### Two types of sparsity for matrices $M \in \mathbb{R}^{n \times p}$ II - Through a factorization of $M = \mathbf{U}\mathbf{V}^{\top}$

- Matrix  $\mathbf{M} = \mathbf{U}\mathbf{V}^{ op}$ ,  $\mathbf{U} \in \mathbb{R}^{n imes k}$  and  $\mathbf{V} \in \mathbb{R}^{p imes k}$
- Low rank: *m* small



 $\bullet$  Sparse decomposition: U sparse



# Structured (sparse) matrix factorizations

• Matrix  $\mathbf{M} = \mathbf{U}\mathbf{V}^{ op}$ ,  $\mathbf{U} \in \mathbb{R}^{n imes k}$  and  $\mathbf{V} \in \mathbb{R}^{p imes k}$ 

#### $\bullet$ Structure on ${\bf U}$ and/or ${\bf V}$

- Low-rank:  ${\bf U}$  and  ${\bf V}$  have few columns
- Dictionary learning / sparse PCA:  ${\bf U}$  has many zeros
- Clustering (k-means):  $\mathbf{U} \in \{0,1\}^{n \times m}$ ,  $\mathbf{U1} = \mathbf{1}$
- Pointwise positivity: non negative matrix factorization (NMF)
- Specific patterns of zeros
- Low-rank + sparse (Candès et al., 2009)

– etc.

### • Many applications

• Many open questions: algorithms, identifiability, evaluation

### Sparse principal component analysis

- Given data  $\mathbf{X} = (\mathbf{x}_1^{\top}, \dots, \mathbf{x}_n^{\top}) \in \mathbb{R}^{p \times n}$ , two views of PCA:
  - Analysis view: find the projection  $\mathbf{d} \in \mathbb{R}^p$  of maximum variance (with deflation to obtain more components)
  - Synthesis view: find the basis  $d_1, \ldots, d_k$  such that all  $x_i$  have low reconstruction error when decomposed on this basis
- For regular PCA, the two views are equivalent



### Sparse principal component analysis

- Given data  $\mathbf{X} = (\mathbf{x}_1^{\top}, \dots, \mathbf{x}_n^{\top}) \in \mathbb{R}^{p \times n}$ , two views of PCA:
  - Analysis view: find the projection  $\mathbf{d} \in \mathbb{R}^p$  of maximum variance (with deflation to obtain more components)
  - Synthesis view: find the basis  $d_1, \ldots, d_k$  such that all  $x_i$  have low reconstruction error when decomposed on this basis
- For regular PCA, the two views are equivalent
- Sparse (and/or non-negative) extensions
  - Interpretability
  - High-dimensional inference
  - Two views are differents
  - For analysis view, see d'Aspremont, Bach, and El Ghaoui (2008)

## Sparse principal component analysis Synthesis view

• Find  $\mathbf{d}_1, \ldots, \mathbf{d}_k \in \mathbb{R}^p$  sparse so that

$$\sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_{i} \in \mathbb{R}^{m}} \left\| \mathbf{x}_{i} - \sum_{j=1}^{k} (\boldsymbol{\alpha}_{i})_{j} \mathbf{d}_{j} \right\|_{2}^{2} = \sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_{i} \in \mathbb{R}^{m}} \left\| \mathbf{x}_{i} - \mathbf{D} \boldsymbol{\alpha}_{i} \right\|_{2}^{2} \text{ is small}$$

- Look for  $\mathbf{A} = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n) \in \mathbb{R}^{k \times n}$  and  $\mathbf{D} = (\mathbf{d}_1, \dots, \mathbf{d}_k) \in \mathbb{R}^{p \times k}$ such that  $\mathbf{D}$  is sparse and  $\|\mathbf{X} - \mathbf{DA}\|_F^2$  is small

# Sparse principal component analysis Synthesis view

• Find  $\mathbf{d}_1, \dots, \mathbf{d}_k \in \mathbb{R}^p$  sparse so that

$$\sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_{i} \in \mathbb{R}^{m}} \left\| \mathbf{x}_{i} - \sum_{j=1}^{k} (\boldsymbol{\alpha}_{i})_{j} \mathbf{d}_{j} \right\|_{2}^{2} = \sum_{i=1}^{n} \min_{\boldsymbol{\alpha}_{i} \in \mathbb{R}^{m}} \left\| \mathbf{x}_{i} - \mathbf{D} \boldsymbol{\alpha}_{i} \right\|_{2}^{2} \text{ is small}$$

- Look for  $\mathbf{A} = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n) \in \mathbb{R}^{k \times n}$  and  $\mathbf{D} = (\mathbf{d}_1, \dots, \mathbf{d}_k) \in \mathbb{R}^{p \times k}$ such that  $\mathbf{D}$  is sparse and  $\|\mathbf{X} - \mathbf{DA}\|_F^2$  is small
- Sparse formulation (Witten et al., 2009; Bach et al., 2008)
  - Penalize/constrain  $\mathbf{d}_j$  by the  $\ell_1$ -norm for sparsity
  - Penalize/constrain  $lpha_i$  by the  $\ell_2$ -norm to avoid trivial solutions

$$\min_{\mathbf{D},\mathbf{A}} \sum_{i=1}^{n} \|\mathbf{x}_{i} - \mathbf{D}\boldsymbol{\alpha}_{i}\|_{2}^{2} + \lambda \sum_{j=1}^{k} \|\mathbf{d}_{j}\|_{1} \text{ s.t. } \forall i, \|\boldsymbol{\alpha}_{i}\|_{2} \leq 1$$

#### **Sparse PCA vs. dictionary learning**

• Sparse PCA:  $\mathbf{x}_i \approx \mathbf{D} \boldsymbol{lpha}_i$ ,  $\mathbf{D}$  sparse



#### **Sparse PCA vs. dictionary learning**

• Sparse PCA:  $\mathbf{x}_i pprox \mathbf{D} \boldsymbol{lpha}_i$ ,  $\mathbf{D}$  sparse



• Dictionary learning:  $\mathbf{x}_i pprox \mathbf{D} \boldsymbol{lpha}_i$ ,  $\boldsymbol{lpha}_i$  sparse





### Structured matrix factorizations (Bach et al., 2008)

$$\min_{\mathbf{D},\mathbf{A}} \sum_{i=1}^{n} \|\mathbf{x}_{i} - \mathbf{D}\boldsymbol{\alpha}_{i}\|_{2}^{2} + \lambda \sum_{j=1}^{k} \|\mathbf{d}_{j}\|_{\star} \text{ s.t. } \forall i, \|\boldsymbol{\alpha}_{i}\|_{\bullet} \leqslant 1$$
$$\min_{\mathbf{D},\mathbf{A}} \sum_{i=1}^{n} \|\mathbf{x}_{i} - \mathbf{D}\boldsymbol{\alpha}_{i}\|_{2}^{2} + \lambda \sum_{i=1}^{n} \|\boldsymbol{\alpha}_{i}\|_{\bullet} \text{ s.t. } \forall j, \|\mathbf{d}_{j}\|_{\star} \leqslant 1$$

- Optimization by alternating minimization (non-convex)
- $\alpha_i$  decomposition coefficients (or "code"),  $\mathbf{d}_j$  dictionary elements
- Two related/equivalent problems:
  - Sparse PCA = sparse dictionary ( $\ell_1$ -norm on  $\mathbf{d}_j$ )
  - Dictionary learning = sparse decompositions ( $\ell_1$ -norm on  $\alpha_i$ ) (Olshausen and Field, 1997; Elad and Aharon, 2006; Lee et al., 2007)

# **Dictionary learning for image denoising**





original image measurements noise

### **Dictionary learning for image denoising**

- Solving the denoising problem (Elad and Aharon, 2006)
  - Extract all overlapping  $8 \times 8$  patches  $\mathbf{x}_i \in \mathbb{R}^{64}$
  - Form the matrix  $\mathbf{X} = (\mathbf{x}_1^{\top}, \dots, \mathbf{x}_n^{\top}) \in \mathbb{R}^{n \times 64}$
  - Solve a matrix factorization problem:

$$\min_{\mathbf{D},\mathbf{A}} ||\mathbf{X} - \mathbf{D}\mathbf{A}||_F^2 = \min_{\mathbf{D},\mathbf{A}} \sum_{i=1}^n ||\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i||_2^2$$

where  ${\bf A}$  is  ${\color{black}{\text{sparse}}}, {\color{black}{\text{and}}} \ {\bf D}$  is the  ${\color{black}{\text{dictionary}}}$ 

- Each patch is decomposed into  $\mathbf{x}_i = \mathbf{D} oldsymbol{lpha}_i$
- Average the reconstruction  $\mathbf{D}\alpha_i$  of each patch  $\mathbf{x}_i$  to reconstruct a full-sized image
- The number of patches n is large (= number of pixels)

#### **Online optimization for dictionary learning**

$$\min_{\mathbf{A}\in\mathbb{R}^{k\times n},\mathbf{D}\in\mathcal{D}}\sum_{i=1}^{n}\left\{||\mathbf{x}_{i}-\mathbf{D}\boldsymbol{\alpha}_{i}||_{2}^{2}+\lambda||\boldsymbol{\alpha}_{i}||_{1}\right\}$$
$$\mathcal{D}\stackrel{\Delta}{=}\{\mathbf{D}\in\mathbb{R}^{p\times k} \text{ s.t. } \forall j=1,\ldots,k, \ ||\mathbf{d}_{j}||_{2}\leqslant1\}.$$

- $\bullet$  Classical optimization alternates between  ${\bf D}$  and  ${\bf A}.$
- Good results, but very slow !

### **Online optimization for dictionary learning**

$$\min_{\mathbf{D}\in\mathcal{D}}\sum_{i=1}^{n}\min_{\boldsymbol{\alpha}_{i}}\left\{||\mathbf{x}_{i}-\mathbf{D}\boldsymbol{\alpha}_{i}||_{2}^{2}+\lambda||\boldsymbol{\alpha}_{i}||_{1}\right\}$$
$$\mathcal{D}\stackrel{\Delta}{=}\{\mathbf{D}\in\mathbb{R}^{p\times k} \text{ s.t. } \forall j=1,\ldots,k, \ ||\mathbf{d}_{j}||_{2}\leqslant1\}.$$

- $\bullet$  Classical optimization alternates between  ${\bf D}$  and  ${\bf A}.$
- Good results, but very slow !
- Online learning (Mairal, Bach, Ponce, and Sapiro, 2009a) can
  - handle potentially infinite datasets
  - adapt to dynamic training sets
  - online code (http://www.di.ens.fr/willow/SPAMS/)

# Denoising result (Mairal, Bach, Ponce, Sapiro, and Zisserman, 2009c)



# Denoising result (Mairal, Bach, Ponce, Sapiro, and Zisserman, 2009c)



#### What does the dictionary D look like?



THE SALINAS VALLEY is in Northern California. It is a long narrow swale between two ranges of mountains, and the Salinas River winds and twists up the center until it falls at last into Monterey Bay.

I remember my childhood names for grasses and secret flowers. I remember where a toad may live and what time the birds awaken in the summer-and what trees and seasons smelled like-how people looked and walked and smelled even. The memory of odors is very rich.

I remember that the Gabilan Mountains to the east of the valley were light gay mountains full of sun and loveliness and a kind of invitation, so that you wanted to climb into their warm foothills almost as you want to climb into the lap of a beloved mother. They were beckoning mountains with a brown grass love. The Santa Lucias stood up against the sky to the west and kept the valley from the open sea, and they were dark and brooding-unfriendly and dangerous. I always found in myself a dread of west and a love of east. Where I ever got such an idea I cannot say, unless it could be that the morning came over the peaks of the Gabilans and the night drifted back from the ridges of the Santa Lucias. It may be that the birth and death of the day had some part in my feeling about the two ranges of mountains.

From both sides of the valley little streams slipped out of the hill canyons and fell into the bed of the Salinas River. In the winter of wet years the streams ran full-freshet, and they swelled the river until sometimes it raged and boiled, bank full, and then it was a destroyer. The river tore the edges of the farm lands and washed whole acres down; it toppled barns and houses into itself, to go floating and bobbing away. It trapped cows and ome pools would be left in deep swirl places under and willows straightened up with the ood de in their upper branches. T sun drove it undergroun Filt was ummer ot a f ted about it how dangerous it was in a we ad and so ter and boast about anything if it's a ou have. Maybe the less you have, the more you are require mer. You to boast

The floor of the Salinas Valley, between the ranges and below the foothills, is fevel because this valley used to be the bottom of a hundred-mile infet from the sea. The river mouth at Moss Landing was centuries ago the entrance to this long inland water. Once, fifty miles down the valley, my father bored a well. The drill came up first with topsoil and then with gravel and then with white sea sand full of shells and even pi...







# Alternative usages of dictionary learning Computer vision

- Use the "code"  $\alpha$  as representation of observations for subsequent processing (Raina et al., 2007; Yang et al., 2009)
- Adapt dictionary elements to specific tasks (Mairal, Bach, Ponce, Sapiro, and Zisserman, 2009b)
  - Discriminative training for weakly supervised pixel classification (Mairal, Bach, Ponce, Sapiro, and Zisserman, 2008)



# Structured sparse methods for matrix factorization Outline

- Learning problems on matrices
- Sparse methods for matrices
  - Sparse principal component analysis
  - Dictionary learning
- Structured sparse PCA
  - Sparsity-inducing norms and overlapping groups
  - Structure on dictionary elements
  - Structure on decomposition coefficients

### **Sparsity-inducing norms**



- Regularizing by a sparsity-inducing norm  $\psi$
- Most popular choice for  $\psi$ 
  - $\ell_1$ -norm:  $\|\boldsymbol{\alpha}\|_1 = \sum_{j=1}^p |\boldsymbol{\alpha}_j|$
  - Lasso (Tibshirani, 1996), basis pursuit (Chen et al., 2001)
  - $\ell_1$ -norm only encodes cardinality
- Structured sparsity
  - Certain patterns are favored
  - Improvement of interpretability and prediction performance

### **Sparsity-inducing norms**

- Another popular choice for  $\psi$ :
  - The  $\ell_1$ - $\ell_2$  norm,

$$\sum_{G \in \mathbf{F}} \|\boldsymbol{\alpha}_G\|_2 = \sum_{G \in \mathbf{F}} \left(\sum_{j \in G} \boldsymbol{\alpha}_j^2\right)^{1/2}, \text{ with } \mathbf{F} \text{ a partition of } \{1, \dots, p\}$$

- The  $\ell_1$ - $\ell_2$  norm sets to zero groups of non-overlapping variables (as opposed to single variables for the  $\ell_1$  -norm)
- For the square loss, group Lasso (Yuan and Lin, 2006)

### **Sparsity-inducing norms**

- Another popular choice for  $\psi$ :
  - The  $\ell_1$ - $\ell_2$  norm,

$$\sum_{G \in \mathbf{F}} \|\boldsymbol{\alpha}_G\|_2 = \sum_{G \in \mathbf{F}} \left(\sum_{j \in G} \boldsymbol{\alpha}_j^2\right)^{1/2}, \text{ with } \mathbf{F} \text{ a partition of } \{1, \dots, p\}$$

- The  $\ell_1$ - $\ell_2$  norm sets to zero groups of non-overlapping variables (as opposed to single variables for the  $\ell_1$ -norm)
- For the square loss, group Lasso (Yuan and Lin, 2006)
- However, the  $\ell_1$ - $\ell_2$  norm encodes **fixed/static prior information**, requires to know in advance how to group the variables
- $\bullet$  What happens if the set of groups  ${\bf F}$  is not a partition anymore?

# Structured Sparsity (Jenatton, Audibert, and Bach, 2009a)

• When penalizing by the  $\ell_1$ - $\ell_2$  norm,

$$\sum_{G \in \mathbf{F}} \|\boldsymbol{\alpha}_G\|_2 = \sum_{G \in \mathbf{F}} \left(\sum_{j \in G} \boldsymbol{\alpha}_j^2\right)^{1/2}$$

- The  $\ell_1$  norm induces sparsity at the group level: \* Some  $\alpha_G$ 's are set to zero
- Inside the groups, the  $\ell_2$  norm does not promote sparsity

### Examples of set of groups ${\bf F}$

• Selection of contiguous patterns on a sequence, p=6



- ${\bf F}$  is the set of blue groups
- Any union of blue groups set to zero leads to the selection of a contiguous pattern
# Structured Sparsity (Jenatton, Audibert, and Bach, 2009a)

• When penalizing by the  $\ell_1$ - $\ell_2$  norm,

$$\sum_{G \in \mathbf{F}} \|\boldsymbol{\alpha}_G\|_2 = \sum_{G \in \mathbf{F}} \left(\sum_{j \in G} \boldsymbol{\alpha}_j^2\right)^{1/2}$$

- The  $\ell_1$  norm induces sparsity at the group level: \* Some  $\alpha_G$ 's are set to zero
- Inside the groups, the  $\ell_2$  norm does not promote sparsity
- $\bullet$  Intuitively, the zero pattern of w is given by

$$\{j \in \{1, \dots, p\}; \ \boldsymbol{\alpha}_j = 0\} = \bigcup_{G \in \mathbf{F}'} G$$
 for some  $\mathbf{F}' \subseteq \mathbf{F}$ 

This intuition is actually true and can be formalized

## Examples of set of groups ${\bf F}$

 $\bullet$  Selection of rectangles on a 2-D grids, p=25



- ${f F}$  is the set of blue/green groups (with their not displayed complements)
- Any union of blue/green groups set to zero leads to the selection of a rectangle

## Examples of set of groups ${\bf F}$

• Selection of diamond-shaped patterns on a 2-D grids, p = 25.



 It is possible to extend such settings to 3-D space, or more complex topologies

# Relationship between F and Zero Patterns (Jenatton, Audibert, and Bach, 2009a)

#### • $\mathbf{F} \rightarrow \text{Zero patterns}$ :

– by generating the union-closure of  ${\bf F}$ 

#### • Zero patterns $\rightarrow$ **F**:

- Design groups  ${\bf F}$  from any  $union\mathchar`-closed\ set$  of  $zero\ \mbox{patterns}$
- Design groups  ${\bf F}$  from any  $intersection\mathchar`-closed\ set\ of\ non\mathchar`-zero\ patterns$

#### **Related work on structured sparsity**

- Specific hierarchical structure (Zhao et al., 2009; Bach, 2008)
- Union-closed (as opposed to intersection-closed) family of nonzero patterns (Jacob, Obozinski, and Vert, 2009)
- Nonconvex penalties based on information-theoretic criteria with greedy optimization (Baraniuk et al., 2008; Huang et al., 2009)
- Link with submodular functions (Bach, 2010)
  - Acting on supports or level sets

# Sparse structured PCA (Jenatton, Obozinski, and Bach, 2009b)

• Learning sparse and structured dictionary elements:

$$\min_{\mathbf{A} \in \mathbb{R}^{k \times n} \\ \mathbf{D} \in \mathbb{R}^{p \times k}} \sum_{i=1}^{n} \|\mathbf{x}_{i} - \mathbf{D}\boldsymbol{\alpha}_{i}\|_{2}^{2} + \lambda \sum_{j=1}^{p} \psi(\mathbf{d}_{j}) \text{ s.t. } \forall i, \|\boldsymbol{\alpha}_{i}\|_{2} \leq 1$$

- Structure of the dictionary elements determined by the choice of overlapping groups  ${f F}$  (and thus  $\psi$ )
- Efficient learning procedures through " $\eta$ -tricks"

- Reweighted 
$$\ell_2$$
:  $\sum_{G \in \mathbf{F}} \|\mathbf{y}_G\|_2 = \min_{\eta_G \ge 0, G \in \mathbf{F}} \frac{1}{2} \sum_{G \in \mathbf{F}} \left\{ \frac{\|\mathbf{y}_G\|_2^2}{\eta_G} + \eta_G \right\}$ 



• NMF obtains partially local features



(unstructured) sparse PCA Structured sparse PCA

 $\bullet$  Enforce selection of convex nonzero patterns  $\Rightarrow$  robustness to occlusion



(unstructured) sparse PCA Structured sparse PCA

 $\bullet$  Enforce selection of convex nonzero patterns  $\Rightarrow$  robustness to occlusion

• Quantitative performance evaluation on classification task



## Dictionary learning vs. sparse structured PCA Exchange roles of D and A

• Sparse structured PCA (sparse and structured dictionary elements):

$$\min_{\mathbf{A}\in\mathbb{R}^{k\times n}\atop \mathbf{D}\in\mathbb{R}^{p\times k}}\sum_{i=1}^{n}\|\mathbf{x}_{i}-\mathbf{D}\boldsymbol{\alpha}_{i}\|_{2}^{2}+\lambda\sum_{j=1}^{k}\psi(\mathbf{d}_{j}) \text{ s.t. }\forall i, \|\boldsymbol{\alpha}_{i}\|_{2} \leq 1.$$

• Dictionary learning with structured sparsity for  $\alpha$ :

$$\min_{\substack{\mathbf{A}\in\mathbb{R}^{k\times n}\\\mathbf{D}\in\mathbb{R}^{p\times k}}}\sum_{i=1}^{n}\|\mathbf{x}_{i}-\mathbf{D}\boldsymbol{\alpha}_{i}\|_{2}^{2}+\lambda\psi(\boldsymbol{\alpha}_{i}) \text{ s.t. } \forall j, \|\mathbf{d}_{j}\|_{2} \leq 1.$$

## Hierarchical dictionary learning (Jenatton, Mairal, Obozinski, and Bach, 2010)

- Structure on codes  $\alpha$  (not on dictionary D)
- Hierarchical penalization:  $\psi(\alpha) = \sum_{G \in \mathbf{F}} \|\alpha_G\|_2$  where groups G in  $\mathbf{F}$  are equal to set of descendants of some nodes in a tree



• Variable selected after its ancestors (Zhao et al., 2009; Bach, 2008)

## Hierarchical dictionary learning Efficient optimization

$$\min_{\substack{\mathbf{A}\in\mathbb{R}^{k\times n}\\\mathbf{D}\in\mathbb{R}^{p\times k}}}\sum_{i=1}^{n}\|\mathbf{x}_{i}-\mathbf{D}\boldsymbol{\alpha}_{i}\|_{2}^{2}+\lambda\psi(\boldsymbol{\alpha}_{i}) \text{ s.t. } \forall j, \|\mathbf{d}_{j}\|_{2} \leq 1.$$

- Minimization with respect to  $oldsymbol{lpha}_i$  : regularized least-squares
  - Many algorithms dedicated to the  $\ell_1\text{-norm}~\psi({\boldsymbol{\alpha}}) = \|{\boldsymbol{\alpha}}\|_1$
- Proximal methods : first-order methods with optimal convergence rate (Nesterov, 2007; Beck and Teboulle, 2009)
  - Requires solving many times  $\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} \boldsymbol{\alpha}\|_2^2 + \lambda \psi(\boldsymbol{\alpha})$
- Tree-structured regularization : Efficient linear time algorithm based on primal-dual decomposition (Jenatton et al., 2010)

# Hierarchical dictionary learning Application to image denoising

- Reconstruction of 100,000  $8\times8$  natural images patches
  - Remove randomly subsampled pixels
  - Reconstruct with matrix factorization and structured sparsity





#### **Application to image denoising - Dictionary tree**



# Hierarchical dictionary learning Modelling of text corpora

- Each document is modelled through word counts
- Low-rank matrix factorization of word-document matrix
- Probabilistic topic models (Blei et al., 2003)
  - Similar structures based on non parametric Bayesian methods (Blei et al., 2004)
  - Can we achieve similar performance with simple matrix factorization formulation?

# Hierarchical dictionary learning Modelling of text corpora

- Each document is modelled through word counts
- Low-rank matrix factorization of word-document matrix
- Probabilistic topic models (Blei et al., 2003)
  - Similar structures based on non parametric Bayesian methods (Blei et al., 2004)
  - Can we achieve similar performance with simple matrix factorization formulation?
- Experiments:
  - Qualitative: NIPS abstracts (1714 documents, 8274 words)
  - Quantitative: newsgroup articles (1425 documents, 13312 words)

#### **Modelling of text corpora - Dictionary tree**



#### Modelling of text corpora

• Comparison on predicting newsgroup article subjects:



#### **Topic models, NMF and matrix factorization**

- Three different views on the same problem
  - Interesting parallels to be made
  - Common problems to be solved
- Structure on dictionary/decomposition coefficients with adapted priors, e.g., nested Chinese restaurant processes (Blei et al., 2004)
- Learning hyperparameters from data
- Identifiability and interpretation/evaluation of results
- Discriminative tasks (Blei and McAuliffe, 2008; Lacoste-Julien et al., 2008; Mairal et al., 2009b)
- Optimization and local minima

## Conclusion

#### • Structured matrix factorization has many applications

- Machine learning
- Image/signal processing, audio/music (Lefèvre et al., 2011)
- Extensions to other tasks

## Application to background subtraction (Mairal, Jenatton, Obozinski, and Bach, 2010)

Background

 $\ell_1$ -norm

Structured norm



## **Ongoing Work - Digital Zooming**



# **Digital Zooming (Couzinie-Devy et al., 2010)**



## **Digital Zooming (Couzinie-Devy et al., 2010)**



# **Digital Zooming (Couzinie-Devy et al., 2010)**



## **Ongoing Work** - Task-driven dictionaries inverse half-toning (Mairal et al., 2010)



# **Ongoing Work** - Task-driven dictionaries inverse half-toning (Mairal et al., 2010)











## Conclusion

- Structured matrix factorization has many applications
  - Machine learning
  - Image/signal processing, audio/music (Lefèvre et al., 2011)
  - Extensions to other tasks

#### • Algorithmic issues

- Large datasets
- Structured sparsity and convex optimization
- Link with submodular functions (Bach, 2010)

#### • Theoretical issues

- Identifiability of structures and features
- Improved predictive performance
- Other approaches to sparsity and structure

#### References

- Y. Amit, M. Fink, N. Srebro, and S. Ullman. Uncovering shared structures in multiclass classification. In *Proceedings of the 24th international conference on Machine Learning (ICML)*, 2007.
- F. Bach. Exploring large feature spaces with hierarchical multiple kernel learning. In Advances in Neural Information Processing Systems, 2008.
- F. Bach. Structured sparsity-inducing norms through submodular functions. In NIPS, 2010.
- F. Bach, J. Mairal, and J. Ponce. Convex sparse matrix factorizations. Technical Report 0812.1869, ArXiv, 2008.
- R. G. Baraniuk, V. Cevher, M. F. Duarte, and C. Hegde. Model-based compressive sensing. Technical report, arXiv:0808.3572, 2008.
- 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.
- L. Benaroya, F. Bimbot, and R. Gribonval. Audio source separation with a single sensor. *IEEE Transactions on Speech and Audio Processing*, 14(1):191, 2006.
- D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. *The Journal of Machine Learning Research*, 3:993–1022, January 2003.
- D. Blei, T.L. Griffiths, M.I. Jordan, and J.B. Tenenbaum. Hierarchical topic models and the nested Chinese restaurant process. *Advances in neural information processing systems*, 16:106, 2004.
- D.M. Blei and J. McAuliffe. Supervised topic models. In *Advances in Neural Information Processing* Systems (NIPS), volume 20, 2008.

- E.J. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *Arxiv preprint* arXiv:0912.3599, 2009.
- S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Review*, 43(1):129–159, 2001.
- A. d'Aspremont, F. Bach, and L. El Ghaoui. Optimal solutions for sparse principal component analysis. *Journal of Machine Learning Research*, 9:1269–1294, 2008.
- M. Elad and M. Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image Processing*, 15(12):3736–3745, 2006.
- C. Févotte, N. Bertin, and J.-L. Durrieu. Nonnegative matrix factorization with the itakura-saito divergence. with application to music analysis. *Neural Computation*, 21(3), 2009.
- J. Huang, T. Zhang, and D. Metaxas. Learning with structured sparsity. In *Proceedings of the 26th 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 26th International Conference on Machine Learning (ICML)*, 2009.
- R. Jenatton, J.Y. Audibert, and F. Bach. Structured variable selection with sparsity-inducing norms. Technical report, arXiv:0904.3523, 2009a.
- R. Jenatton, G. Obozinski, and F. Bach. Structured sparse principal component analysis. Technical report, arXiv:0909.1440, 2009b.
- R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for sparse hierarchical dictionary learning. In *Submitted to ICML*, 2010.
- S. Lacoste-Julien, F. Sha, and M.I. Jordan. DiscLDA: Discriminative learning for dimensionality

reduction and classification. Advances in Neural Information Processing Systems (NIPS) 21, 2008.

- H. Lee, A. Battle, R. Raina, and A. Ng. Efficient sparse coding algorithms. In Advances in Neural Information Processing Systems (NIPS), 2007.
- J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Discriminative learned dictionaries for local image analysis. In *IEEE Conference on Computer Vision and Pattern Recognition, (CVPR)*, 2008.
- J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online dictionary learning for sparse coding. In *International Conference on Machine Learning (ICML)*, 2009a.
- J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Supervised dictionary learning. *Advances in Neural Information Processing Systems (NIPS)*, 21, 2009b.
- J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Non-local sparse models for image restoration. In *International Conference on Computer Vision (ICCV)*, 2009c.
- J. Mairal, R. Jenatton, G. Obozinski, and F. Bach. Network flow algorithms for structured sparsity. In *NIPS*, 2010.
- 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.
- M. Pontil, A. Argyriou, and T. Evgeniou. Multi-task feature learning. In *Advances in Neural Information Processing Systems*, 2007.
- R. Raina, A. Battle, H. Lee, B. Packer, and A.Y. Ng. Self-taught learning: Transfer learning from unlabeled data. In *Proceedings of the 24th International Conference on Machine Learning (ICML)*, 2007.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society. Series B*, pages 267–288, 1996.
- D.M. Witten, R. Tibshirani, and T. Hastie. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534, 2009.
- J. Yang, K. Yu, Y. Gong, and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- 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, G. Rocha, and B. Yu. Grouped and hierarchical model selection through composite absolute penalties. *Annals of Statistics*, 37(6A):3468–3497, 2009.