

# Representation learning: Dictionary learning

## Formation ENSTA-ParisTech

### Conférence IA

Florent Chatelain \*    Olivier Michel \*

\* Univ. Grenoble Alpes, GIPSA-lab

10-13 February 2020

# Outline

Dictionary problem

Sparse coding

K-SVD algorithm

Dictionary learning examples

## Reminder of regression and sparsity

- ▶  $\mathbf{x} \in \mathbb{R}^p \leftarrow$  observations,
- ▶  $D \in \mathbb{R}^{p \times K} \leftarrow$  design/regression matrix,

### Regression problem

Find coefficient vector  $\mathbf{a} \in \mathbb{R}^K$  s.t.  $\mathbf{x} \approx D\mathbf{a}$ , i.e.  $\|\mathbf{x} - D\mathbf{a}\| \leq \epsilon$

🔍 regularization/model selection: **sparsity** constraint

$$\min_{\mathbf{a}} \|\mathbf{x} - D\mathbf{a}\|_2^2 \quad \text{s.t. } \text{Pen}(\mathbf{a}) \leq T, \quad \text{with e.g.}$$

- ▶  $\text{Pen}(\mathbf{a}) = \|\mathbf{a}\|_1 \leftarrow \ell_1$  norm,
- ▶  $\text{Pen}(\mathbf{a}) = \|\mathbf{a}\|_0 \equiv \#\{i : 1 \leq i \leq K \text{ and } \mathbf{a}_i \neq 0\} \leftarrow \ell_0$  pseudo-norm, i.e. number of non-zero components for  $\mathbf{a}$

## Dictionary problem

### Assumptions

- Inputs:  $\mathbf{x}_i \in \mathbb{R}^P$ , for  $i = 1, \dots, n$ ,  $\leftarrow$  learning set

### Objective

“Inverse” the regression problem, i.e. find  $D \in \mathbb{R}^{P \times K}$  ensuring a **sparse** approximation/representation of the learning set  $\mathbf{x}_1, \dots, \mathbf{x}_n$

$$\mathbf{x}_i \approx D\mathbf{a}_i, \text{ for } i = 1, \dots, n, \text{ s.t.}$$

- $K \ll n \leftarrow$  sparsity of the representation 'basis' ( $\sim$  low rank constraint)
- $\|\mathbf{a}_i\|_0 \leq T$ , for  $i = 1, \dots, n \leftarrow$  sparsity of the coefficients

### Dictionary terminology

- $D \equiv (\mathbf{d}_1 | \dots | \mathbf{d}_K) \in \mathbb{R}^{P \times K} \leftarrow$  **dictionary** s.t.  $K \ll n$  (low rank)
- $\mathbf{d}_k \in \mathbb{R}^P \leftarrow$  **atoms** with  $k = 1, \dots, K$  (column vectors of  $D$ )
- $\mathbf{a}_i \in \mathbb{R}^K \leftarrow$  **coefficients** of  $\mathbf{x}_i$  s.t.  $\|\mathbf{a}_i\|_0 \leq T$  (sparsity)

## Dictionary vs Orthonormal Basis (ONB)

### Orthonormal Bases (ONB)

Lots of good properties for ONB (e.g. Fourier, orthonormal wavelets, ...):

- ▶ uniqueness of the decomposition ( $n$  coefficients to represent a  $n$ -dimensional vector),
- ▶ existence of fast transforms,
- ▶ projections to obtain approximate representations, ...

**But** not necessarily good for sparse representations

- 👉 e.g. sums of sinusoids + spikes (+ steps) are not sparse in either Fourier or identity bases...

### Dictionary

Change of mindset

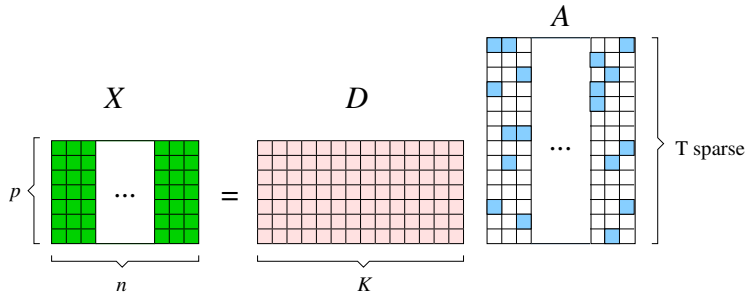
- ▶ loose uniqueness of representation using **overcomplete/redundant** bases
- 👉 get **sparser** representations in return!

## Recap: Dictionary and Sparse representations

### Hypotheses

► Inputs:  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $1 \leq i \leq n$ ,  $\rightarrow X = (\mathbf{x}_1 | \dots | \mathbf{x}_n) \in \mathbb{R}^{p \times n}$

Objective: sparse representation/factorization  $X \approx DA$



►  $D \equiv (\mathbf{d}_1 | \dots | \mathbf{d}_K) \in \mathbb{R}^{p \times K} \leftarrow$  dictionary (redundant), with  $K \ll n$ .

►  $\mathbf{d}_k \in \mathbb{R}^p \leftarrow$  atoms with  $k = 1, \dots, K$

►  $A \equiv (\mathbf{a}_1 | \dots | \mathbf{a}_n) \in \mathbb{R}^{K \times n} \leftarrow$  coefficients s.t.  $\mathbf{x}_i \approx D\mathbf{a}_i$  with  $\|\mathbf{a}_i\|_0 \leq T$

👁 only the observations  $\mathbf{x}_i$  are known: unsupervised problem

## Formalizing the dictionary learning problem

### Cost criterion with sparsity constraints

$$\min_{D, \{\mathbf{a}_i\}} J(D, \{\mathbf{a}_i\}) \equiv \sum_{i=1}^n \|\mathbf{x}_i - D\mathbf{a}_i\|_2^2 = \|X - DA\|_F^2,$$

s.t.  $\text{Pen}(\mathbf{a}_i) \leq T$  for  $i = 1, \dots, n$

- ▶  $X = (\mathbf{x}_1 | \dots | \mathbf{x}_n) \in \mathbb{R}^{p \times n}$ ,  $D = (\mathbf{d}_1 | \dots | \mathbf{d}_K) \in \mathbb{R}^{p \times K}$ ,  $A = (\mathbf{a}_1 | \dots | \mathbf{a}_n) \in \mathbb{R}^{K \times n}$ ,
- ▶  $\|M\|_F^2 \equiv \sum_{i,j} m_{ij}^2 \leftarrow$  Frobenius norm

### Identifiability issue

- Pb:** atoms  $\mathbf{d}_k$  and their associate coefficients  $\mathbf{a}^k \in \mathbb{R}^{1 \times n}$  ( $k$ th line of  $\mathbf{a}$ ) are defined up to a scaling factor in the factorization criterion
- 👉 atoms are assumed to be **normalized**, i.e.  $\|\mathbf{d}_k\|_2 = 1$ , for  $k = 1, \dots, K$

## Optimizing the factorization criterion

### Alternate minimization method

$$\min_{D, \{\mathbf{a}_i\}} J(D, \{\mathbf{a}_i\}) = \|X - DA\|_F^2, \quad \text{s.t. } \text{Pen}(\mathbf{a}_i) \leq T$$

Joint minimization w.r.t.  $A$  and  $D$  split into 2 simpler separated steps

1. **sparse coding** to estimate the coefficients  $A$  for a given dictionary  $D$ ,
  2. **dictionary  $D$  update** performed atom-by-atom, for a given  $A$
- 👉 repeat 1) and 2) until a stopping criterion

### Optimization issues

- ▶ when  $\text{Pen}(\mathbf{a}_i) = \|\mathbf{a}_i\|_0 \leftarrow$  **NP-hard** combinatorics problem...
- ▶ joint minimization w.r.t.  $D$  and  $\{\mathbf{a}_i\}_{1 \leq i \leq n} \leftarrow$  **non-convex** problem
- 👉 **sub-optimal** solution around an initial value for  $D$  (cf clustering)



## Outline

Dictionary problem

Sparse coding

K-SVD algorithm

Dictionary learning examples

## Sparse coding: update of the coefficients $\mathbf{a}$

**Assumption:**  $D$  fixed

**Objective:** find sparse  $A^*$  minimizing  $J(A) \equiv J(D, A) = \|X - DA\|_F^2$

👉 linear regression problem on each  $\mathbf{a}_i$  with sparsity constraints

### Basis Pursuit (BP)

Convex relaxation of the  $\ell_0$  pseudo-norm with the  $\ell_1$  norm:

$$A^* = \min_A \|X - DA\|_F^2 + \lambda \sum_{i=1}^n \|\mathbf{a}_i\|, \quad \text{where } \lambda > 0,$$

► convex problem on  $A$ , with Lasso solutions for  $\mathbf{a}_i$ ,  $i = 1, \dots, n$

👉 Online Dictionary Learning (ODL) algorithm (Mairal *et al.*, 2009)

### (Orthogonal) Matching Pursuit (OMP)

Sparse approximate (sub-optimal) solution of the  $\ell_0$  minimization problem using a greedy algorithm to select the “best matching” atoms

## Orthogonal Matching Pursuit

Let  $D_\Gamma$  be the sub dictionary composed of the only atoms  $\mathbf{d}_k$  s.t.  $k \in \Gamma$

1. **Initialization:**  $\Gamma = \emptyset$ ,  $\mathbf{r} = \mathbf{x}_i$  (sweep on the observations  $\mathbf{x}_i$ )
2. **for**  $iter = 1, \dots, T$  **do**
3.   Select the best matching atom, i.e. most correlated with the residuals

$$\hat{k} \leftarrow \max_{1 \leq k \leq K} |\mathbf{r}^T \mathbf{d}_k|$$

4.   Update the active set:  $\Gamma \leftarrow \Gamma \cup \{\hat{k}\}$
5.   Update the residuals:  $\mathbf{r} \leftarrow \left( I_p - D_\Gamma (D_\Gamma^T D_\Gamma)^{-1} D_\Gamma^T \right) \mathbf{x}_i$ ,
6. **endfor**
7. **Output:**  $\mathbf{a}_i \equiv$  coefficients of the orthogonal projection of  $\mathbf{x}_i$  on the space spanned by  $D_\Gamma$

☞ residuals **orthogonal** to the sparse signal approximation  $D_\Gamma \mathbf{a}_i$

☞ by construction,  $\text{card}(\Gamma) = T \Rightarrow \|\mathbf{a}_i\|_0 \leq T$

## Outline

Dictionary problem

Sparse coding

**K-SVD algorithm**

Dictionary learning examples

## K-SVD algorithm outline

1. preliminary tools: SVD and low rank approximation
2. K-SVD dictionary update and K-SVD algorithm

## Singular Value Decomposition (SVD)

Let  $X \in \mathbb{R}^{n \times m}$  be a real rectangular matrix. There exists a factorization, called a singular value decomposition of  $X$ , of the form

$$X = U \Sigma V^T, \quad \text{where}$$

- ▶  $U \in \mathbb{R}^{n \times n}$  is **orthonormal** ( $UU^T = U^T U = I_n$ )  $\leftarrow$  matrix of left-singular vectors  $\mathbf{u}_k \in \mathbb{R}^n$  s.t.  
 $U = (\mathbf{u}_1 | \dots | \mathbf{u}_n)$ ,
  - ▶  $V \in \mathbb{R}^{m \times m}$  is **orthonormal** ( $VV^T = V^T V = I_m$ )  $\leftarrow$  matrix of right-singular vectors  $\mathbf{v}_k \in \mathbb{R}^m$  s.t.  
 $V = (\mathbf{v}_1 | \dots | \mathbf{v}_m)$ ,
  - ▶  $\Sigma \in \mathbb{R}^{n \times m}$  is a rectangular **diagonal** matrix with non negative diagonal entries  $\Sigma_{ii} \equiv \lambda_i \geq 0$  for  $i = 1, \dots, \min(n, m)$
- 👉 the  $\lambda_i$ 's are **uniquely defined** and called the **singular values** of  $X$
- 👉 By convention, these singular values are sorted in descending order:

$$\lambda_1 \geq \dots \geq \lambda_{\min(n,m)} \geq 0$$

## Singular value decomposition (SVD) and principal components analysis (PCA)

$X \in \mathbb{R}^{n \times m}$  with SVD  $X = U\Sigma V^T$

### Eigendecomposition of Gram matrices

Gram matrices express as

$$\begin{aligned} X^T X &= V \Sigma^T U^T U \Sigma V^T = V (\Sigma^T \Sigma) V^T, \\ XX^T &= U \Sigma V^T V \Sigma^T U^T = U (\Sigma \Sigma^T) U^T, \end{aligned}$$

where  $\Sigma^T \Sigma \in \mathbb{R}^{m \times m}$  and  $\Sigma \Sigma^T \in \mathbb{R}^{n \times n}$  are square diagonal

- ▶ right-singular vectors of  $X$  are **eigenvectors** of  $X^T X = V (\Sigma^T \Sigma) V^T$ ,
  - ▶ left-singular vectors of  $X$  are **eigenvectors** of  $XX^T = U (\Sigma \Sigma^T) U^T$ ,
  - ▶ non-zero singular values  $\lambda_i$ 's are the **square roots of the non-zero eigenvalues** of  $X^T X$  or  $XX^T$
- 👉 SVD yields **principal components analysis (PCA)** decomposition

## Singular Value Decomposition (SVD) and low rank approximation

- ▶  $X \in \mathbb{R}^{n \times m}$  with SVD  $X = U\Sigma V^T$
- 👉 left-singular vectors  $\mathbf{u}_k \in \mathbb{R}^n$  of  $X$  are the columns of  $U$ , for  $k = 1, \dots, n$
- 👉 right-singular vectors  $\mathbf{v}_k \in \mathbb{R}^m$  of  $X$  are the columns of  $V$ , for  $k = 1, \dots, m$

### Eckart–Young–Mirsky theorem

For the Frobenius norm ( $\|M\|_F^2 = \sum_{i,j} m_{i,j}^2$ ), the solution to the **low-rank approximation** problem

$\min_{\hat{X}} \|X - \hat{X}\|_F^2$  s.t.  $\text{rank}(\hat{X}) \leq r$ , is

$$\hat{X} = \sum_{k=1}^r \lambda_k \mathbf{u}_k \mathbf{v}_k^T,$$

where  $\lambda_k$ ,  $\mathbf{u}_k$  and  $\mathbf{v}_k$  are the first **singular values** and **left/right-singular vectors** of  $X$ , for  $k = 1, \dots, r \leq \min(m, n)$ .



## Dictionary update: K-SVD algorithm (Aharon *et al.*, 2006)

**Objective:** For a given  $A$ , find  $D$  minimizing  $J(D) \equiv J(D, A) = \|X - DA\|_F^2$

Sweep on the atoms  $\mathbf{d}_k$  that are sequentially updated, for  $k = 1, \dots, K$ , as

$$\min_{\mathbf{d}_k} \|X - DA\|_F^2 = \|E^k - \mathbf{d}_k \mathbf{a}^k\|_F^2$$

- ▶  $\mathbf{a}^k \in \mathbb{R}^{1 \times n}$  is the  $k$ th **line** of  $A$ ,
- ▶  $E^k = X - \sum_{l \neq k} \mathbf{d}_l \mathbf{a}^l \in \mathbb{R}^{p \times n}$ ,
- ▶  $\mathbf{d}_k \mathbf{a}^k \in \mathbb{R}^{p \times n}$  is a **rank 1 matrix**

### Atoms (and coefficients) update

Best rank 1 approximation of  $E^k$  obtained from its **SVD** as  $\lambda_1 \mathbf{u}_1 \mathbf{v}_1^T$ ,

- ▶  $\lambda_1 \equiv$  largest singular value,  $\mathbf{u}_1$  and  $\mathbf{v}_1^T \equiv$  associated singular vectors
- 👉  $\mathbf{d}_k = \mathbf{u}_1 \leftarrow$  unit norm vector by construction
- 👉  $\mathbf{a}^k = \lambda_1 \mathbf{v}_1^T \leftarrow$  **Pb:** there is no reason for the updated  $\mathbf{a}^k$  to be **sparse!**

## Dictionary update: K-SVD algorithm (Aharon *et al.*, 2006)

### K-SVD solution to enforce sparsity of the coefficients

For updating the  $k$ th atom restrict attention to the only observations  $\mathbf{x}_i$  for which the atom is active, i.e. data and dictionary columns with indexes in

$$\Gamma_k = \{i : \alpha_i^k \neq 0\}$$

- ▶  $\mathbf{a}_{\Gamma_k}^k \in \mathbb{R}^{|\Gamma_k|}$  is the vector of the only non zero components of  $\mathbf{a}^k \in \mathbb{R}^n$
- ▶  $E_{\Gamma_k}^k \in \mathbb{R}^{n \times |\Gamma_k|}$  is the matrix obtained by retaining the only column vectors of  $E^k = X - \sum_{l \neq k} \mathbf{d}_l \mathbf{a}^l$  with column index  $i \in \Gamma_k$
- ▶  $\lambda_1 \equiv$  largest singular value,  $\mathbf{u}_1$  and  $\mathbf{v}_1^T \equiv$  associated singular vectors of the SVD of  $E_{\Gamma_k}^k$

### K-SVD atoms (and coefficients) update

- ▶  $\mathbf{d}_k \leftarrow \mathbf{u}_1$  (unit norm vector by construction)
- ▶  $\mathbf{a}_{\Gamma_k}^k \leftarrow \lambda_1 \mathbf{v}_1^T$  (hence  $\mathbf{a}^k$  remains as **sparse** as sparse coding outputs)
- 🔊 both atoms and coefficients are updated during the K-SVD dictionary update step

## K-SVD algorithm

0. **Initialize** dictionary  $D^{(0)}$  using union of ONBs elements / randomly by picking  $K$  observations
1. **Sparse coding**: computation of the coefficients  $\mathbf{a}_i$  for each  $\mathbf{x}_i$  given  $D^{(t-1)}$  (usually using OMP with  $\|\mathbf{a}_i\|_0 \leq T$ )
2. **Dictionary update**: for each atom  $\mathbf{d}_k \in D^{(t-1)}$ ,  $1 \leq k \leq K$ 
  - ▶ Define the set of active entries  $\Gamma_k = \{i : \alpha_i^k \neq 0\}$
  - ▶ Compute  $E^k = X - \sum_{l \neq k} \mathbf{d}_l \mathbf{a}_l^T$  and its restriction  $E_{\Gamma_k}^k$  to columns in  $\Gamma_k$
  - ▶ Compute the first largest singular value  $\lambda_1 > 0$  of  $E_{\Gamma_k}^k$ , and the associate left and right singular vectors  $\mathbf{u}_1$  and  $\mathbf{v}_1$
  - ▶ Update the  $k$ -th atom and the associated coefficients
 
$$\begin{array}{ll} \mathbf{d}_k & \leftarrow \mathbf{u}_1, \\ \mathbf{a}_{\Gamma_k}^k & \leftarrow \lambda_1 \mathbf{v}_1^T \end{array}$$
- ▶ Repeat steps 1 and 2 for  $t = 1, \dots$  until convergence

## Properties of K-SVD algorithm

### Convergence

Convergence to a **local minimum** guaranteed if the sparse coding step always decrease the empirical quadratic error  $\|X - DA\|_F^2$

- ▶ no theoretical guaranties when using approximations like OMP
- ▶ in practice, condition empirically satisfied when  $T \ll n$

### Generalization of $K$ -means

- ▶ For  $T = 1$  sparsity constraints, normalizing the coefficients  $\mathbf{a}_i$  rather than the atoms  $\mathbf{d}_k$ , then  $K$ -SVD reduces to the  $K$ -means algorithm
  - 💡 dictionary atoms are the cluster means/centroids
- ▶ In the general case  $T \geq 1$ , the computation of the mean is replaced by a SVD computation for updating the atoms
  - 💡 “ $K$ -SVD” algorithm

## Implementation details of K-SVD

### Sparse coding

- ▶ OMP is preferred for efficiency,
- ▶ for denoising problem, more convenient to solve (approximatively) the similar problem

$$\min_{\mathbf{a}_i} \|\mathbf{a}_i\|_0 \quad \text{s.t.} \quad \|\mathbf{x}_i - D\mathbf{a}_i\| \leq \epsilon,$$

where  $\epsilon$  can be tuned from the noise power  $\leftarrow$  change of stopping criterion in OMP

### Dictionary update heuristics

- ▶ Pruning atoms that are not “used” enough,
- ▶ Removing atoms too coherent with each other,
- ✎ pruned or removed atom replaced by the least well-explained observation, i.e. observation  $\mathbf{x}_l$  where  $l = \arg \max_{1 \leq i \leq n} \|\mathbf{x}_i - D\mathbf{a}_i\|_2$

## Outline

Dictionary problem

Sparse coding

K-SVD algorithm

Dictionary learning examples

## Examples of sparse representations for image restoration

### Solving the denoising problem [Elad and Aharon, 2006]

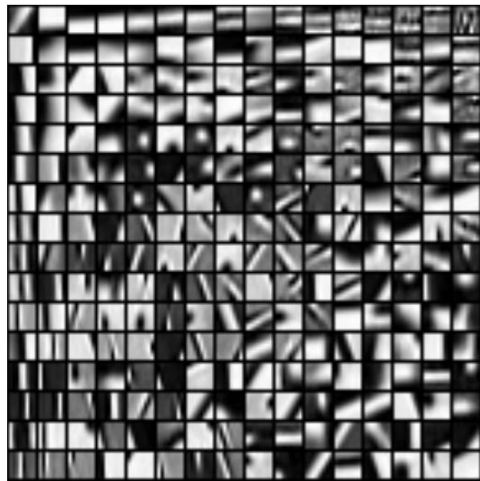
- ▶ Extract all overlapping  $8 \times 8$  patches  $\mathbf{x}_i$ , for  $1 \leq i \leq n$  with  $n > 10^5$
- ▶ Solve a matrix factorization (dictionary learning) problem:

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

with **sparsity** constraints  $\text{Pen}(\mathbf{a}_i) \leq T$ , e.g.  $\text{Pen}(\mathbf{a}_i) = \|\mathbf{a}_i\|_0$

- ▶ Average the reconstruction of each patch.

## K-SVD results for image restoration



- Dictionary trained on a **noisy** version of the image

From ICCV tutorial

[http://lear.inrialpes.fr/people/mairal/tutorial\\_iccv09/tuto\\_part2.pdf](http://lear.inrialpes.fr/people/mairal/tutorial_iccv09/tuto_part2.pdf)



## Applications: denoising

[Mairal, Bach, Ponce, Sapiro, and Zisserman, 2009c]



- Dictionary trained on a **noisy** version of the image

From ICCV tutorial

[http://lear.inrialpes.fr/people/mairal/tutorial\\_iccv09/tuto\\_part2.pdf](http://lear.inrialpes.fr/people/mairal/tutorial_iccv09/tuto_part2.pdf)

## Applications: inpainting (1)

Inpainting, [Mairal, Elad, and Sapiro, 2008a]



- Dictionary trained on the image with **missing** data

From ICCV tutorial [http://lear.inrialpes.fr/people/mairal/tutorial\\_iccv09/tutorial\\_part2.pdf](http://lear.inrialpes.fr/people/mairal/tutorial_iccv09/tutorial_part2.pdf)

## Applications: inpainting (2)

Inpainting, [Mairal, Elad, and Sapiro, 2008a]



- Dictionary trained on the image with **missing** data

From ICCV tutorial [http://lear.inrialpes.fr/people/mairal/tutorial\\_iccv09/tuto\\_part2.pdf](http://lear.inrialpes.fr/people/mairal/tutorial_iccv09/tuto_part2.pdf)