



HAL
open science

Sparse Modeling for Image and Vision Processing

Julien Mairal, Francis Bach, Jean Ponce

► **To cite this version:**

Julien Mairal, Francis Bach, Jean Ponce. Sparse Modeling for Image and Vision Processing. now publishers, 8 (2-3), pp.85-283, 2014, Foundations and Trends in Computer Graphics and Vision, 978-1-68083-008-8. 10.1561/06000000058 . hal-01081139v2

HAL Id: hal-01081139

<https://inria.hal.science/hal-01081139v2>

Submitted on 6 Dec 2014

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Sparse Modeling for Image and Vision Processing

Julien Mairal
Inria¹
julien.mairal@inria.fr

Francis Bach
Inria²
francis.bach@inria.fr

Jean Ponce
Ecole Normale Supérieure³
jean.ponce@ens.fr

¹LEAR team, laboratoire Jean Kuntzmann, CNRS, Univ. Grenoble Alpes, France.

²SIERRA team, département d'informatique de l'Ecole Normale Supérieure, ENS/CNRS/Inria UMR 8548, France.

³WILLOW team, département d'informatique de l'Ecole Normale Supérieure, ENS/CNRS/Inria UMR 8548, France.

Contents

1	A Short Introduction to Parsimony	2
1.1	Early concepts of parsimony in statistics	6
1.2	Wavelets in signal processing	8
1.3	Modern parsimony: the ℓ_1 -norm and other variants	13
1.4	Dictionary learning	32
1.5	Compressed sensing and sparse recovery	36
1.6	Theoretical results about dictionary learning	40
2	Discovering the Structure of Natural Images	45
2.1	Pre-processing	47
2.2	Principal component analysis	54
2.3	Clustering or vector quantization	57
2.4	Dictionary learning	61
2.5	Structured dictionary learning	62
2.6	Other matrix factorization methods	66
2.7	Discussion	74
3	Sparse Models for Image Processing	76
3.1	Image denoising	77
3.2	Image inpainting	84
3.3	Image demosaicking	86

3.4	Image up-scaling	90
3.5	Inverting nonlinear local transformations	93
3.6	Video processing	97
3.7	Face compression	97
3.8	Other patch modeling approaches	102
4	Sparse Coding for Visual Recognition	107
4.1	A coding and pooling approach to image modeling	108
4.2	The botany of sparse feature coding	116
4.3	Face recognition	123
4.4	Patch classification and edge detection	125
4.5	Connections with neural networks	132
4.6	Other applications	136
5	Optimization Algorithms	142
5.1	Sparse reconstruction with the ℓ_0 -penalty	143
5.2	Sparse reconstruction with the ℓ_1 -norm	150
5.3	Iterative reweighted- ℓ_1 methods	156
5.4	Iterative reweighted- ℓ_2 methods	159
5.5	Optimization for dictionary learning	161
5.6	Other optimization techniques	172
6	Conclusions	173
	Acknowledgments	175
	References	176

Abstract

In recent years, a large amount of multi-disciplinary research has been conducted on sparse models and their applications. In statistics and machine learning, the sparsity principle is used to perform model selection—that is, automatically selecting a simple model among a large collection of them. In signal processing, sparse coding consists of representing data with linear combinations of a few dictionary elements. Subsequently, the corresponding tools have been widely adopted by several scientific communities such as neuroscience, bioinformatics, or computer vision. The goal of this monograph is to offer a self-contained view of sparse modeling for visual recognition and image processing. More specifically, we focus on applications where the dictionary is learned and adapted to data, yielding a compact representation that has been successful in various contexts.

1

A Short Introduction to Parsimony

In its most general definition, the principle of sparsity, or parsimony, consists of representing some phenomenon with as few variables as possible. It appears to be central to many research fields and is often considered to be inspired from an early doctrine formulated by the philosopher and theologian William of Ockham in the 14th century, which essentially favors simple theories over more complex ones. Of course, the link between Ockham and the tools presented in this monograph is rather thin, and more modern views seem to appear later in the beginning of the 20th century. Discussing the scientific method, Wrinch and Jeffreys [1921] introduce indeed a simplicity principle related to parsimony as follows:

The existence of simple laws is, then, apparently, to be regarded as a quality of nature; and accordingly we may infer that it is justifiable to prefer a simple law to a more complex one that fits our observations slightly better.

Remarkably, Wrinch and Jeffreys [1921] further discuss statistical modeling of physical observations and relate the concept of “simplicity” to the number of learning parameters; as a matter of fact, this concept is relatively close to the contemporary view of parsimony.

Subsequently, numerous tools have been developed by statisticians to build models of physical phenomena with good predictive power. Models are usually learned from observed data, and their generalization performance is evaluated on test data. Among a collection of plausible models, the simplest one is often preferred, and the number of underlying parameters is used as a criterion to perform model selection [Mallows, 1964, 1966, Akaike, 1973, Hocking, 1976, Barron et al., 1998, Rissanen, 1978, Schwarz, 1978, Tibshirani, 1996].

In signal processing, similar problems as in statistics arise, but a different terminology is used. Observations, or data vectors, are called “signals”, and data modeling appears to be a crucial step for performing various operations such as restoration, compression, or for solving inverse problems. Here also, the sparsity principle plays an important role and has been successful [Mallat and Zhang, 1993, Pati et al., 1993, Donoho and Johnstone, 1994, Cotter et al., 1999, Chen et al., 1999, Donoho, 2006, Candès et al., 2006]. Each signal is approximated by a sparse linear combination of prototypes called dictionary elements, resulting in simple and compact models.

However, statistics and signal processing remain two distinct fields with different objectives and methodology; specifically, signals often come from the same data source, *e.g.*, natural images, whereas problems considered in statistics are unrelated to each other in general. Then, a long series of works has been devoted to finding appropriate dictionaries for signal classes of interest, leading to various sorts of wavelets [Freeman and Adelson, 1991, Simoncelli et al., 1992, Donoho, 1999, Candès and Donoho, 2002, Do and Vetterli, 2005, Le Pennec and Mallat, 2005, Mallat, 2008]. Even though statistics and signal processing have devised most of the methodology of sparse modeling, the parsimony principle was also discovered independently in other fields. To some extent, it appears indeed in the work of Markowitz [1952] about portfolio selection in finance, and also in geophysics [Claerbout and Muir, 1973, Taylor et al., 1979].

In neuroscience, Olshausen and Field [1996, 1997] proposed a significantly different approach to sparse modeling than previously established practices. Whereas classical techniques in signal

processing were using fixed off-the-shelf dictionaries, the method of Olshausen and Field [1996, 1997] consists of learning it from training data. In a pioneer exploratory experiment, they demonstrated that dictionary learning could easily discover underlying structures in natural image patches; later, their approach found numerous applications in many fields, notably in image and audio processing [Lewicki, 2002, Elad and Aharon, 2006, Mairal et al., 2009, Yang et al., 2010a] and computer vision [Raina et al., 2007, Yang et al., 2009, Zeiler et al., 2011, Mairal et al., 2012, Song et al., 2012, Castrodad and Sapiro, 2012, Elhamifar et al., 2012, Pokrass et al., 2013].

The goal of this monograph is to present basic tools of sparse modeling and their applications to visual recognition and image processing. We aim at offering a self-contained view combining pluri-disciplinary methodology, practical advice, and a large review of the literature. Most of the figures in the paper are produced with the software SPAMS¹, and the corresponding Matlab code will be provided on the first author’s webpage.

The monograph is organized as follows: the current introductory section is divided into several parts providing a simple historical view of sparse estimation. In Section 1.1, we start with early concepts of parsimony in statistics and information theory from the 70’s and 80’s. We present the use of sparse estimation within the wavelet framework in Section 1.2, which was essentially developed in the 90’s. Section 1.3 introduces the era of “modern parsimony”—that is, the ℓ_1 -norm and its variants, which have been heavily used during the last two decades. Section 1.4 is devoted to the dictionary learning formulation originally introduced by Olshausen and Field [1996, 1997], which is a key component of most applications presented later in this monograph. In Sections 1.5 and 1.6, we conclude our introductory tour with some theoretical aspects, such as the concept of “compressed sensing” and sparse recovery that has attracted a large attention in recent years.

With all these parsimonious tools in hand, we discuss the use of sparse coding and related sparse matrix factorization techniques for discovering the underlying structure of natural image patches in Sec-

¹available here <http://spams-devel.gforge.inria.fr/>.

tion 2 . Even though the task here is subjective and exploratory, it is the first successful instance of dictionary learning; the insight gained from these early experiments forms the basis of concrete applications presented in subsequent sections.

Section 3 covers numerous applications of sparse models of natural image patches in image processing, such as image denoising, super-resolution, inpainting, or demosaicking. This section is concluded with other related patch-modeling approaches.

Section 4 presents recent success of sparse models for visual recognition, such as codebook learning of visual descriptors, face recognition, or more low-level tasks such as edge detection and classification of textures and digits. We conclude the section with other computer vision applications such as visual tracking and data visualization.

Section 5 is devoted to optimization algorithms. It presents in a concise way efficient algorithms for solving sparse decomposition and dictionary learning problems.

We see our monograph as a good complement of other books and monographs about sparse estimation, which offer different perspectives, such as Mallat [2008], Elad [2010] in signal and image processing, or Bach et al. [2012a] in optimization and machine learning. We also put the emphasis on the structure of natural image patches learned with dictionary learning, and thus present an alternative view to the book of Hyvärinen et al. [2009], which is focused on independent component analysis.

Notation. In this monograph, vectors are denoted by bold lower-case letters and matrices by upper-case ones. For instance, we consider in the rest of this paragraph a vector \mathbf{x} in \mathbb{R}^n and a matrix \mathbf{X} in $\mathbb{R}^{m \times n}$. The columns of \mathbf{X} are represented by indexed vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ such that we can write $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$. The i -th entry of \mathbf{x} is denoted by $\mathbf{x}[i]$, and the i -th entry of the j -th column of \mathbf{X} is represented by $\mathbf{X}[i, j]$. For any subset g of $\{1, \dots, n\}$, we denote by $\mathbf{x}[g]$ the vector in $\mathbb{R}^{|g|}$ that records the entries of \mathbf{x} corresponding to indices in g . For $q \geq 1$, we define the ℓ_q -norm of \mathbf{x} as $\|\mathbf{x}\|_q \triangleq (\sum_{i=1}^n |\mathbf{x}[i]|^q)^{1/q}$, and the ℓ_∞ -norm as $\|\mathbf{x}\|_\infty \triangleq \lim_{q \rightarrow +\infty} \|\mathbf{x}\|_q = \max_{i=1, \dots, n} |\mathbf{x}[i]|$.

For $q < 1$, we define the ℓ_q -penalty as $\|\mathbf{x}\|_q \triangleq \sum_{i=1}^n |\mathbf{x}[i]|^q$, which, with an abuse of terminology, is often referred to as ℓ_q -norm. The ℓ_0 -penalty simply counts the number of non-zero entries in a vector: $\|\mathbf{x}\|_0 \triangleq \#\{i \text{ s.t. } \mathbf{x}[i] \neq 0\}$. For a matrix \mathbf{X} , we define the Frobenius norm $\|\mathbf{X}\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n \mathbf{X}[i, j]^2\right)^{1/2}$. When dealing with a random variable X defined on a probability space, we denote its expectation by $\mathbb{E}[X]$, assuming that there is no measurability or integrability issue.

1.1 Early concepts of parsimony in statistics

A large number of statistical procedures can be formulated as maximum likelihood estimation. Given a statistical model with parameters $\boldsymbol{\theta}$, it consists of minimizing with respect to $\boldsymbol{\theta}$ an objective function representing the negative log-likelihood of observed data. Assuming for instance that we observe independent samples $\mathbf{z}_1, \dots, \mathbf{z}_n$ of the (unknown) data distribution, we need to solve

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \left[\mathcal{L}(\boldsymbol{\theta}) \triangleq - \sum_{i=1}^n \log P_{\boldsymbol{\theta}}(\mathbf{z}_i) \right], \quad (1.1)$$

where P is some probability distribution parameterized by $\boldsymbol{\theta}$.

Simple methods such as ordinary least squares can be written as (1.1). Consider for instance data points \mathbf{z}_i that are pairs (y_i, \mathbf{x}_i) , with y_i is an observation in \mathbb{R} and \mathbf{x}_i is a vector in \mathbb{R}^p , and assume that there exists a linear relation $y_i = \mathbf{x}_i^\top \boldsymbol{\theta} + \varepsilon_i$, where ε_i is an approximation error for observation i . Under a model where the ε_i 's are independent and identically normally distributed with zero-mean, Eq. (1.1) is equivalent to a least square problem:²

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \sum_{i=1}^n \frac{1}{2} \left(y_i - \mathbf{x}_i^\top \boldsymbol{\theta} \right)^2.$$

To prevent overfitting and to improve the interpretability of the learned model, it was suggested in early work that a solution involving only a

²Note that the Gaussian noise assumption is not necessary to justify the ordinary least square formulation. It is only sufficient to interpret it as maximum likelihood estimation. In fact, as long as the conditional expectation $\mathbb{E}[y|\mathbf{x}]$ is linear, the ordinary least square estimator is statistically consistent under mild assumptions.

few model variables could be more appropriate than an exact solution of (1.1); in other words, a sparse solution involving only—let us say— k variables might be desirable in some situations. Unfortunately, such a strategy yields two difficulties: first, it is not clear a priori how to choose k ; second, finding the best subset of k variables is NP-hard in general [Natarajan, 1995]. The first issue was addressed with several criteria for controlling the trade-off between the sparsity of the solution $\boldsymbol{\theta}$ and the adequacy of the fit to training data. For the second issue, approximate computational techniques have been proposed.

Mallows’s C_p , AIC, and BIC. For the ordinary least squares problem, Mallows [1964, 1966] introduced the C_p -statistics, later generalized by Akaike [1973] with the Akaike information criterion (AIC), and then by Schwarz [1978] with the Bayesian information criterion (BIC). Using C_p , AIC, or BIC is equivalent to solving the penalized ℓ_0 -maximum likelihood estimation problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \mathcal{L}(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_0, \quad (1.2)$$

where λ depends on the chosen criterion [see Hastie et al., 2009], and $\|\boldsymbol{\theta}\|_0$ is the ℓ_0 -penalty. Similar formulations have also been derived by using the minimum description length (MDL) principle for model selection [Rissanen, 1978, Barron et al., 1998]. As shown by Natarajan [1995], the problem (1.2) is NP-hard, and approximate algorithms are necessary unless p is very small, *e.g.*, $p < 30$.

Forward selection and best subset selection for least squares. To obtain an approximate solution of (1.2), a classical approach is the forward selection technique, which is a greedy algorithm that solves a sequence of maximum likelihood estimation problems computed on a subset of variables. After every iteration, a new variable is added to the subset according to the chosen sparsity criterion in a greedy manner. Some variants allow backward steps—that is, a variable can possibly exit the active subset after an iteration. The algorithm is presented in more details in Section 5.1 and seems to be due to Efroymson [1960], according to Hocking [1976]. Other approaches considered in the 70’s include

also the *leaps and bounds* technique of Furnival and Wilson [1974], a branch-and-bound algorithm providing the exact solution of (1.2) with exponential worst-case complexity.

1.2 Wavelets in signal processing

In signal processing, similar problems as in statistics have been studied in the context of wavelets. In a nutshell, a wavelet basis represents a set of functions ϕ_1, ϕ_2, \dots that are essentially dilated and shifted versions of each other. Unlike Fourier basis, wavelets have the interesting properties to be localized both in the space and frequency domains, and to be suitable to multi-resolution analysis of signals [Mallat, 1989].

The concept of parsimony is central to wavelets. When a signal f is “smooth” in a particular sense [see Mallat, 2008], it can be well approximated by a linear combination of a few wavelets. Specifically, f is close to an expansion $\sum_i \alpha_i \phi_i$ where only a few coefficients α_i are non-zero, and the resulting compact representation has effective applications in estimation and compression. The wavelet theory is well developed for continuous signals, *e.g.*, f is chosen in the Hilbert space $L^2(\mathbb{R})$, but also for discrete signals f in \mathbb{R}^m , making it suitable to modern digital image processing.

Since the first wavelet was introduced by Haar [1910], much research has been devoted to designing a wavelet set that is adapted to particular signals such as natural images. After a long quest for finding good orthogonal basis such as the one proposed by Daubechies [1988], a series of works has focused on wavelet sets whose elements are not linearly independent. It resulted a large number of variants, such as steerable wavelets [Simoncelli et al., 1992], curvelets [Candès and Donoho, 2002], contourlets [Do and Vetterli, 2005], or bandlets [Le Pennec and Mallat, 2005]. For the purpose of our monograph, one concept related to sparse estimation is particularly important; it is called *wavelet thresholding*.

Sparse estimation and wavelet thresholding. Let us consider a discrete signal represented by a vector \mathbf{x} in \mathbb{R}^p and an orthogonal wavelet basis set $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ —that is, satisfying $\mathbf{D}^\top \mathbf{D} = \mathbf{I}$ where \mathbf{I} is the

identity matrix. Approximating \mathbf{x} by a sparse linear combination of wavelet elements can be formulated as finding a sparse vector $\boldsymbol{\alpha}$ in \mathbb{R}^p , say with k non-zero coefficients, that minimizes

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_0 \leq k. \quad (1.3)$$

The sparse decomposition problem (1.3) is an instance of the best subset selection formulation presented in Section 1.1 where $\boldsymbol{\alpha}$ represents model parameters, demonstrating that similar topics arise in statistics and signal processing. However, whereas (1.3) is NP-hard for general matrices \mathbf{D} [Natarajan, 1995], we have assumed \mathbf{D} to be orthogonal in the context of wavelets. As such, (1.3) is equivalent to

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{D}^\top \mathbf{x} - \boldsymbol{\alpha}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_0 \leq k,$$

and admits a closed form. Let us indeed define the vector $\boldsymbol{\beta} \triangleq \mathbf{D}^\top \mathbf{x}$ in \mathbb{R}^p , corresponding to the exact non-sparse decomposition of \mathbf{x} onto \mathbf{D} —that is, we have $\mathbf{x} = \mathbf{D}\boldsymbol{\beta}$ since \mathbf{D} is orthogonal. To obtain the best k -sparse approximation, we denote by μ the k -th largest value among the set $\{|\boldsymbol{\beta}[1]|, \dots, |\boldsymbol{\beta}[p]|\}$, and the solution $\boldsymbol{\alpha}^{\text{ht}}$ of (1.3) is obtained by applying to $\boldsymbol{\beta}$ an operator called “hard-thresholding” and defined as

$$\boldsymbol{\alpha}^{\text{ht}}[i] = \mathbf{1}_{|\boldsymbol{\beta}[i]| \geq \mu} \boldsymbol{\beta}[i] = \begin{cases} \boldsymbol{\beta}[i] & \text{if } |\boldsymbol{\beta}[i]| \geq \mu, \\ 0 & \text{otherwise,} \end{cases} \quad (1.4)$$

where $\mathbf{1}_{|\boldsymbol{\beta}[i]| \geq \mu}$ is the indicator function, which is equal to 1 if $|\boldsymbol{\beta}[i]| \geq \mu$ and 0 otherwise. In other words, the hard-thresholding operator simply sets to zero coefficients from $\boldsymbol{\beta}$ whose magnitude is below the threshold μ . The corresponding procedure, called “wavelet thresholding”, is simple and effective for image denoising, even though it does not perform as well as recent state-of-the-art techniques presented in Section 3. When an image \mathbf{x} is noisy, *e.g.*, corrupted by white Gaussian noise, and μ is well chosen, the estimate $\mathbf{D}\boldsymbol{\alpha}^{\text{ht}}$ is a good estimate of the clean original image. The terminology “hard” is defined in contrast to an important variant called the “soft-thresholding operator”, which was introduced by Donoho and Johnstone [1994] in the context of wavelets:³

³Note that the soft-thresholding operator appears in fact earlier in the statistics literature [see Efron and Morris, 1971, Bickel, 1984], but it was used there for a

$$\alpha^{\text{st}}[i] \triangleq \text{sign}(\beta[i]) \max(|\beta[i]| - \lambda, 0) = \begin{cases} \beta[i] - \lambda & \text{if } \beta[i] \geq \lambda, \\ \beta[i] + \lambda & \text{if } \beta[i] \leq -\lambda, \\ 0 & \text{otherwise,} \end{cases} \quad (1.5)$$

where λ is a parameter playing the same role as μ in (1.4). Not only does the operator set small coefficients of β to zero, but it also reduces the magnitude of the non-zero ones. Both operators are illustrated and compared to each other in Figure 1.1. Interestingly, whereas α^{ht} is the solution of (1.3) when μ corresponds to the entry of $\beta = \mathbf{D}^\top \mathbf{x}$ with k -th largest magnitude, α^{st} is in fact the solution of the following sparse reconstruction problem with the orthogonal matrix \mathbf{D} :

$$\min_{\alpha \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\alpha\|_2^2 + \lambda \|\alpha\|_1. \quad (1.6)$$

This formulation will be the topic of the next section for general non-orthogonal matrices. Similar to statistics where choosing the parameter k of the best subset selection was difficult, automatically selecting the best thresholds μ or λ has been a major research topic [see, *e.g.* Donoho and Johnstone, 1994, 1995, Chang et al., 2000a,b].

Structured group thresholding. Wavelets coefficients have a particular structure since the basis elements \mathbf{d}_i are dilated and shifted versions of each other. It is for instance possible to define neighborhood relationships for wavelets whose spatial supports are close to each other, or hierarchical relationships between wavelets with same or similar localization but with different scales. For one-dimensional signals, we present in Figure 1.2 a typical organization of wavelet coefficients on a tree with arrows representing such relations. For two-dimensional images, the structure is slightly more involved and the coefficients are usually organized as a collection of quadrees [see Mallat, 2008, for more details]; we present such a configuration in Figure 1.3.

A natural idea has inspired the recent concept of *group sparsity* that will be presented in the next section; it consists in exploiting the wavelet structure to improve thresholding estimators. Specifically, it is

different purpose.

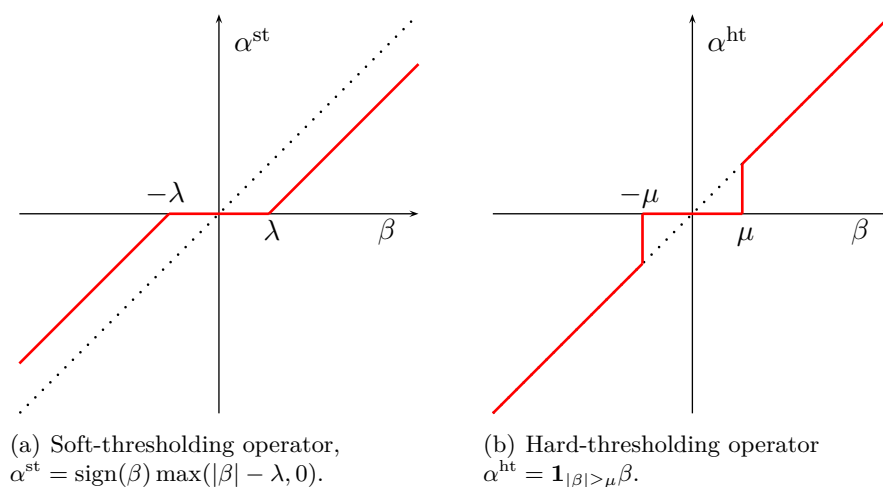


Figure 1.1: Soft- and hard-thresholding operators, which are commonly used for signal estimation with orthogonal wavelet basis.

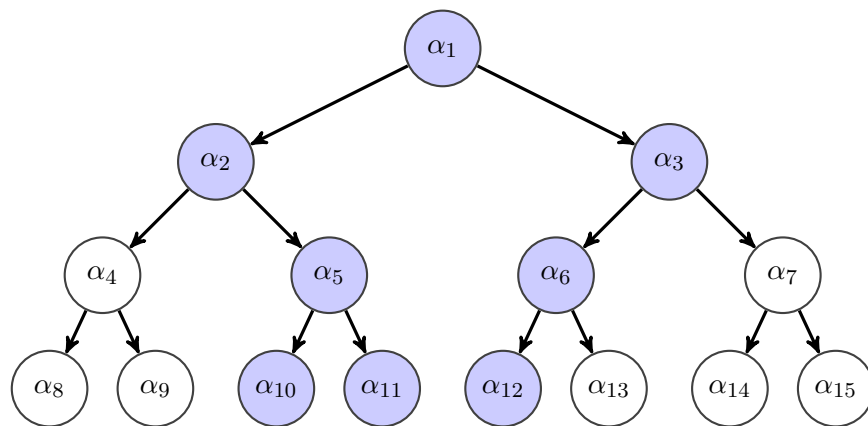


Figure 1.2: Illustration of a wavelet tree with four scales for one-dimensional signals. Nodes represent wavelet coefficients and their depth in the tree correspond to the scale parameter of the wavelet. We also illustrate the zero-tree coding scheme [Shapiro, 1993] in this figure: Empty nodes correspond to zero coefficient: according to the zero-tree coding scheme, their descendants in the tree are also zero.

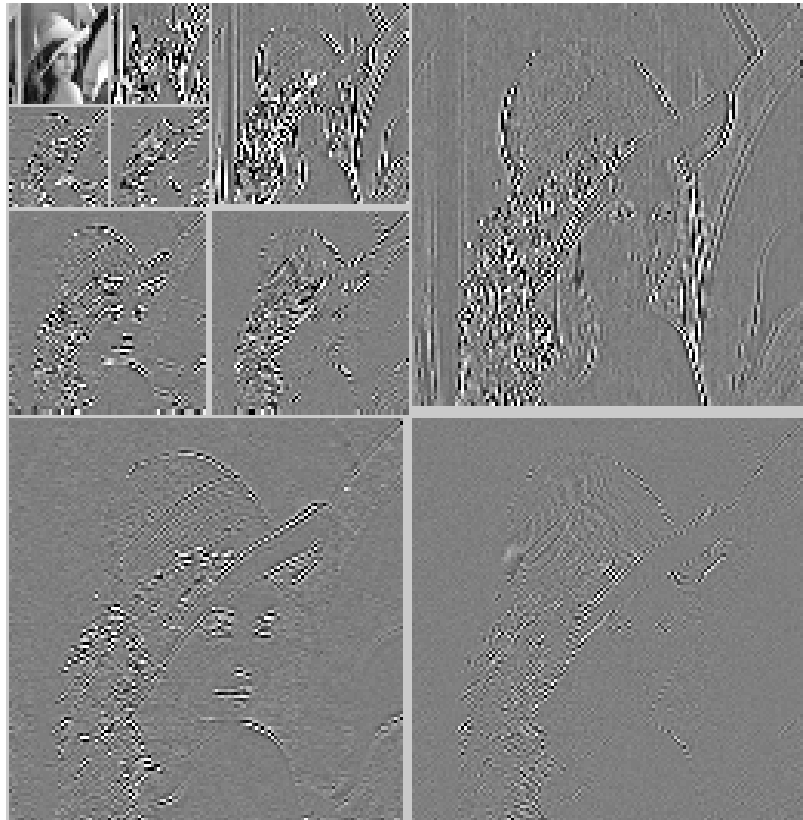


Figure 1.3: Wavelet coefficients displayed for the image *lena* using the orthogonal basis of Daubechies [1988]. A few coefficients representing a low-resolution version of the image are displayed on the top-left corner. Wavelets corresponding to this low-resolution image are obtained by filtering the original image with shifted versions of a low-pass filter called “scaling function” or “father wavelet”. The rest of the coefficients are organized into three quadrees (on the right, on the left, and on the diagonal). Each quadtree is obtained by filtering the original image with a wavelet at three different scales and at different positions. The value zero is represented by the grey color; negative values appear in black, and positive values in white. The wavelet decomposition and this figure have been produced with the software package `matlabPyrTools` developed by Eero Simoncelli and available here: <http://www.cns.nyu.edu/~lcv/software.php>.

possible to use neighborhood relations between wavelet basis elements to define groups of coefficients that form a partition \mathcal{G} of $\{1, \dots, p\}$, and use a group-thresholding operator [Hall et al., 1999, Cai, 1999] defined for every group g in \mathcal{G} as

$$\boldsymbol{\alpha}^{\text{gt}}[g] \triangleq \begin{cases} \left(1 - \frac{\lambda}{\|\boldsymbol{\beta}[g]\|_2}\right) \boldsymbol{\beta}[g] & \text{if } \|\boldsymbol{\beta}[g]\|_2 \geq \lambda, \\ 0 & \text{otherwise,} \end{cases} \quad (1.7)$$

where $\boldsymbol{\beta}[g]$ is the vector of size $|g|$ recording the entries of $\boldsymbol{\beta}$ whose indices are in g . By using such an estimator, groups of neighbor coefficients are set to zero together when their joint ℓ_2 -norm falls below the threshold λ . Interestingly, even though the next interpretation does not appear in early work about group-thresholding [Hall et al., 1999, Cai, 1999], it is possible to view $\boldsymbol{\alpha}^{\text{gt}}$ with $\boldsymbol{\beta} = \mathbf{D}^\top \mathbf{x}$ as the solution of the following penalized problem

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \sum_{g \in \mathcal{G}} \|\boldsymbol{\alpha}[g]\|_2, \quad (1.8)$$

where the closed-form solution (1.7) holds because \mathbf{D} is orthogonal [see Bach et al., 2012a]. Such a formulation will be studied in the next section for general matrices.

Finally, other ideas for exploiting both structure and wavelet parsimony have been proposed. One is a coding scheme called “zero-tree” wavelet coding [Shapiro, 1993], which uses the tree structure of wavelets to force all descendants of zero coefficients to be zero as well. Equivalently, a coefficient can be non-zero only if its parent in the tree is non-zero, as illustrated in Figure 1.2. This idea has been revisited later in a more general context by Zhao et al. [2009]. Other complex models have been used as well for modeling interactions between coefficients: we can mention the application of hidden Markov models (HMM) to wavelets by Crouse et al. [1998] and the Gaussian scale mixture model of Portilla et al. [2003].

1.3 Modern parsimony: the ℓ_1 -norm and other variants

The era of “modern” parsimony corresponds probably to the use of convex optimization techniques for solving feature selection or

sparse decomposition problems. Even though the ℓ_1 -norm was introduced for that purpose in geophysics [Claerbout and Muir, 1973, Taylor et al., 1979], it was popularized in statistics with the Lasso estimator of Tibshirani [1996] and independently in signal processing with the basis pursuit formulation of Chen et al. [1999]. Given observations \mathbf{x} in \mathbb{R}^n and a matrix of predictors \mathbf{D} in $\mathbb{R}^{n \times p}$, the Lasso consists of learning a linear model $\mathbf{x} \approx \mathbf{D}\boldsymbol{\alpha}$ by solving the following quadratic program:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_1 \leq \mu. \quad (1.9)$$

As detailed in the sequel, the ℓ_1 -norm encourages the solution $\boldsymbol{\alpha}$ to be sparse and the parameter μ is used to control the trade-off between data fitting and the sparsity of $\boldsymbol{\alpha}$. In practice, reducing the value of μ leads indeed to sparser solution in general, *i.e.*, with more zeroes, even though there is no formal relation between the sparsity of $\boldsymbol{\alpha}$ and its ℓ_1 -norm for general matrices \mathbf{D} .

The basis pursuit denoising formulation of Chen et al. [1999] is relatively similar but the ℓ_1 -norm is used as a penalty instead of a constraint. It can be written as

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1, \quad (1.10)$$

which is essentially equivalent to (1.9) from a convex optimization perspective, and in fact (1.10) is also often called ‘‘Lasso’’ in the literature. Given some data \mathbf{x} , matrix \mathbf{D} , and parameter $\mu > 0$, we indeed know from Lagrange multiplier theory [see, *e.g.*, Borwein and Lewis, 2006, Boyd and Vandenberghe, 2004] that for all solution $\boldsymbol{\alpha}^*$ of (1.9), there exists a parameter $\lambda \geq 0$ such that $\boldsymbol{\alpha}^*$ is also a solution of (1.10). We note, however, that there is no direct mapping between λ and μ , and thus the choice of formulation (1.9) or (1.10) should be made according to how easy it is to select the parameters λ or μ . For instance, one may prefer (1.9) when a priori information about the ℓ_1 -norm of the solution is available. In Figure 1.4, we illustrate the effect of changing the value of the regularization parameter λ on the solution of (1.10) for two datasets. When $\lambda = 0$, the solution is dense; in general, increasing λ sets more and more variables to zero. However, the relation

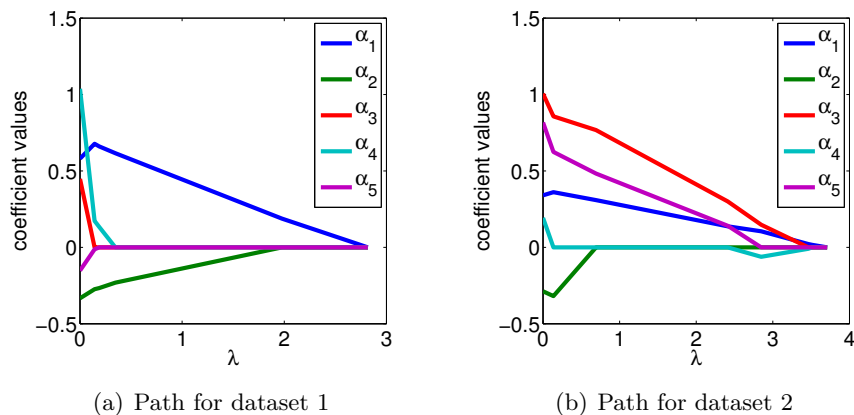


Figure 1.4: Two examples of regularization paths for the Lasso/Basis Pursuit. The curves represent the values of the $p = 5$ entries of the solutions of (1.10) when varying the parameter λ for two datasets. On the left, the relation between λ and the sparsity of the solution is monotonic; On the right, this is not the case. Note that the paths are piecewise linear, see Section 5.2 for more details.

between λ and the sparsity of the solution is not exactly monotonic. In a few cases, increasing λ yields a denser solution.

Another “equivalent” formulation consists of finding a sparse decomposition under a reconstruction constraint:

$$\min_{\alpha \in \mathbb{R}^p} \|\alpha\|_1 \quad \text{s.t.} \quad \|\mathbf{x} - \mathbf{D}\alpha\|_2^2 \leq \varepsilon. \quad (1.11)$$

This formulation can be useful when we have a priori knowledge about the noise level and the parameter ε is easy to choose. The link between (1.10) and (1.11) is similar to the link between (1.10) and (1.9).

For noiseless problems, Chen et al. [1999] have also introduced a formulation simply called “basis pursuit” (without the terminology “denoising”), defined as

$$\min_{\alpha \in \mathbb{R}^p} \|\alpha\|_1 \quad \text{s.t.} \quad \mathbf{x} = \mathbf{D}\alpha, \quad (1.12)$$

which is related to (1.10) in the sense that the set of solutions of (1.10) converges to the solutions of (1.12) when λ converges to 0^+ , whenever the linear system $\mathbf{x} = \mathbf{D}\alpha$ is feasible. These four formulations (1.9-1.12)

have gained a large success beyond the statistics and signal processing communities. More generally, the ℓ_1 -norm has been used as a regularization function beyond the least-square context, leading to problems of the form

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} f(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|_1, \quad (1.13)$$

where $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is a loss function. In the rest of this section, we will present several variants of the ℓ_1 -norm, but before that, we will try to understand why such a penalty is appropriate for sparse estimation.

Why does the ℓ_1 -norm induce sparsity? Even though we have claimed that there is no rigorous relation between the sparsity of $\boldsymbol{\alpha}$ and its ℓ_1 -norm in general, intuition about the sparsity-inducing effect of the ℓ_1 -norm may be obtained from several viewpoints.

Analytical point of view. In the previous section about wavelets, we have seen that when \mathbf{D} is orthogonal, the ℓ_1 -decomposition problem (1.10) admits an analytic closed form solution (1.5) obtained by soft-thresholding. As a result, whenever the magnitude of the inner product $\mathbf{d}_i^\top \mathbf{x}$ is smaller than λ for an index i , the corresponding variable $\boldsymbol{\alpha}^*[i]$ is equal to zero. Thus, the number of zeroes of the solution $\boldsymbol{\alpha}^*$ monotonically increases with λ .

For non-orthogonal matrices \mathbf{D} , such a monotonic relation does not formally hold anymore; in practice, the sparsity-inducing property of the ℓ_1 -penalty remains effective, as illustrated in Figure 1.4. Some intuition about this fact can be gained by studying optimality conditions for the general ℓ_1 -regularized problem (1.13) where f is a differentiable function. The following lemma details these conditions.

Lemma 1.1 (Optimality conditions for ℓ_1 -regularized problems).

A vector $\boldsymbol{\alpha}^*$ in \mathbb{R}^p is a solution of (1.13) if and only if

$$\forall i = 1, \dots, p \quad \begin{cases} -\nabla f(\boldsymbol{\alpha}^*)[i] = \lambda \operatorname{sign}(\boldsymbol{\alpha}^*[i]) & \text{if } \boldsymbol{\alpha}^*[i] \neq 0, \\ |\nabla f(\boldsymbol{\alpha}^*)[i]| \leq \lambda & \text{otherwise.} \end{cases} \quad (1.14)$$

Proof. A proof using the classical concept of subdifferential from convex optimization can be found in [see, e.g., Bach et al., 2012a]. Here,

we provide instead an elementary proof using the simpler concept of directional derivative for nonsmooth functions, defined as, when the limit exists,

$$\nabla g(\boldsymbol{\alpha}, \boldsymbol{\kappa}) \triangleq \lim_{t \rightarrow 0^+} \frac{g(\boldsymbol{\alpha} + t\boldsymbol{\kappa}) - g(\boldsymbol{\alpha})}{t},$$

for a function $g : \mathbb{R}^p \rightarrow \mathbb{R}$ at a point $\boldsymbol{\alpha}$ in \mathbb{R}^p and a direction $\boldsymbol{\kappa}$ in \mathbb{R}^p . For convex functions g , directional derivatives always exist and a classical optimality condition for $\boldsymbol{\alpha}^*$ to be a minimum of g is to have $\nabla g(\boldsymbol{\alpha}^*, \boldsymbol{\kappa})$ non-negative for all directions $\boldsymbol{\kappa}$ [Borwein and Lewis, 2006]. Intuitively, this means that one cannot find any direction $\boldsymbol{\kappa}$ such that an infinitesimal move along $\boldsymbol{\kappa}$ from $\boldsymbol{\alpha}^*$ decreases the value of the objective. When g is differentiable, the condition is equivalent to the classical optimality condition $\nabla g(\boldsymbol{\alpha}^*) = 0$.

We can now apply the directional derivative condition to the function $g : \boldsymbol{\alpha} \mapsto f(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|_1$, which is equivalent to

$$\forall \boldsymbol{\kappa} \in \mathbb{R}^p, \quad \nabla f(\boldsymbol{\alpha}^*)^\top \boldsymbol{\kappa} + \lambda \sum_{i=1}^p \left\{ \begin{array}{ll} \text{sign}(\boldsymbol{\alpha}^*[i])\boldsymbol{\kappa}[i] & \text{if } \boldsymbol{\alpha}^*[i] \neq 0, \\ |\boldsymbol{\kappa}[i]| & \text{otherwise} \end{array} \right\} \geq 0. \quad (1.15)$$

It is then easy to show that (1.15) holds for all $\boldsymbol{\kappa}$ if and only if the inequality holds for the specific values $\boldsymbol{\kappa} = \mathbf{e}_i$ and $\boldsymbol{\kappa} = -\mathbf{e}_i$ for all i , where \mathbf{e}_i is the vector in \mathbb{R}^p with zeroes everywhere except for the i -th entry that is equal to one. This immediately provides an equivalence between (1.15) and (1.14). \square

Lemma 1.1 is interesting from a computational point of view (see Section 5.2), but it also tells us that when $\lambda \geq \|\nabla f(0)\|_\infty$, the conditions (1.14) are satisfied for $\boldsymbol{\alpha}^* = 0$, the sparsest solution possible.

Physical point of view. In image processing or computer vision, the word “energy” often denotes the objective function of a minimization problem; it is indeed common in physics to have complex systems that stabilize at a configuration of minimum potential energy. The negative of the energy’s gradient represents a force, a terminology we will borrow in this paragraph. Consider for instance a one-dimensional ℓ_1 -

regularized estimation problem

$$\min_{\alpha \in \mathbb{R}} \frac{1}{2}(\beta - \alpha)^2 + \lambda|\alpha|, \quad (1.16)$$

where β is a positive constant. Whenever α is non-zero, the ℓ_1 -penalty is differentiable with derivative $\lambda \text{sign}(\alpha)$. When interpreting this objective as an energy minimization problem, the ℓ_1 -penalty can be seen as applying *a force driving α towards the origin with constant intensity λ* . Consider now instead the squared ℓ_2 -penalty, also called regularization of Tikhonov [1963], or ridge regression regularization [Hoerl and Kennard, 1970]:

$$\min_{\alpha \in \mathbb{R}} \frac{1}{2}(\beta - \alpha)^2 + \frac{\lambda}{2}\alpha^2. \quad (1.17)$$

The derivative of the quadratic energy $(\lambda/2)\alpha^2$ is $\lambda\alpha$. It can be interpreted as *a force that also points to the origin but with linear intensity $\lambda|\alpha|$* . Therefore, the force corresponding to the ridge regularization can be arbitrarily strong when α is large, but it fades away when α gets close to zero. As a result, the squared ℓ_2 -regularization does not have a sparsity-inducing effect. From an analytical point of view, we have seen that the solution of (1.16) is zero when $|\beta|$ is smaller than λ . In contrast, the solution of (1.17) admits a closed form $\alpha^* = \beta/(1 + \lambda)$. And thus, regardless of the parameter λ , the solution is never zero.

We present a physical example illustrating this phenomenon in Figure 1.5. We use springs whose potential energy is known to be quadratic, and objects with a gravitational potential energy that is approximately linear on the Earth's surface.

Geometrical point of view. The sparsity-inducing effect of the ℓ_1 -norm can also be interpreted by studying the geometry of the ℓ_1 -ball $\{\alpha \in \mathbb{R}^p : \|\alpha\|_1 \leq \mu\}$. More precisely, understanding the effect of the Euclidean projection onto this set is important: in simple cases where the design matrix \mathbf{D} is orthogonal, the solution of (1.9) can indeed be obtained by the projection

$$\min_{\alpha \in \mathbb{R}^p} \frac{1}{2}\|\beta - \alpha\|_2^2 \quad \text{s.t.} \quad \|\alpha\|_1 \leq \mu, \quad (1.18)$$

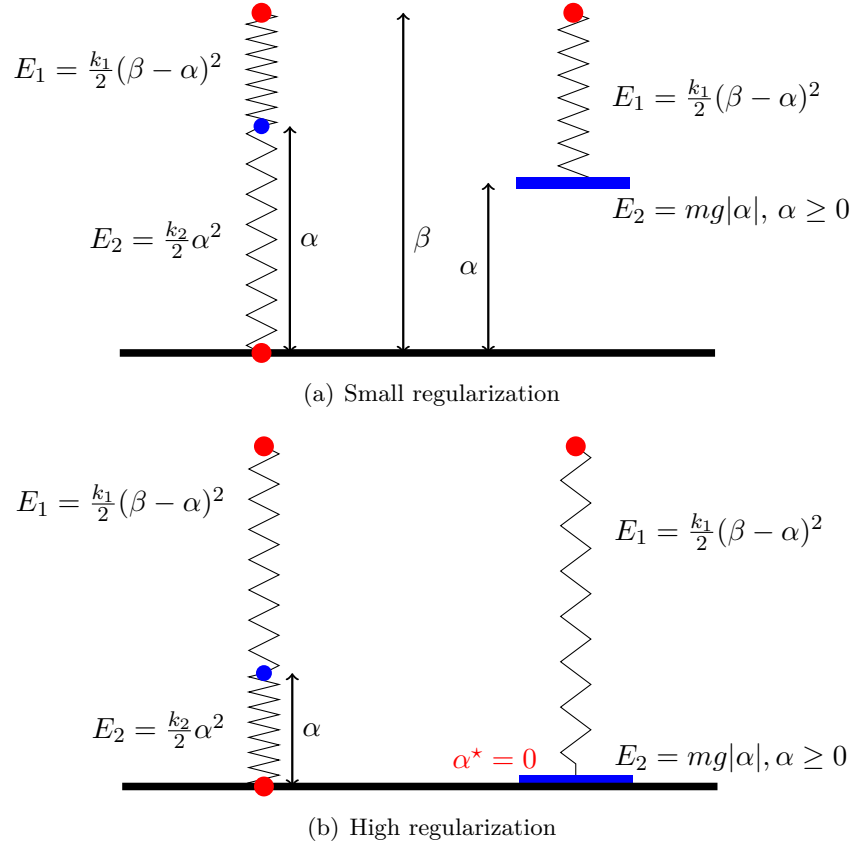


Figure 1.5: A physical system illustrating the sparsity-inducing effect of the ℓ_1 -norm (on the right) in contrast to the Tikhonov-ridge regularization (on the left). Three springs are represented in each figure, two on the left, one on the right. Red points are fixed and cannot move. On the left, two springs are linked to each other by a blue point whose position can vary. On the right, a blue object of mass m is attached to the spring. Right and left configurations define two different dynamical systems with energies $E_1 + E_2$; on the left, E_1 and E_2 are elastic potential energies; on the right, E_1 is the same as on the left, whereas E_2 is a gravitational potential energy, where g is the gravitational constant on the Earth's surface. Both system can evolve according to their initial positions, and stabilize for the value of α^* that minimizes the energy $E_1 + E_2$, assuming that some energy can be dissipated by friction forces. On the left, it is possible to show that $\alpha^* = \beta k_1 / (k_1 + k_2)$ and thus, the solution α^* is never equal to zero, regardless of the strength k_2 of the bottom spring. On the right, the solution is obtained by soft-thresholding: $\alpha^* = \max(\beta - mg/k_1, 0)$. As shown on Figure 1.5(b), when the mass m is large enough, the blue object touches the ground and $\alpha^* = 0$. Figure adapted from [Mairal, 2010].

where $\beta = \mathbf{D}^\top \mathbf{x}$. When \mathbf{D} is not orthogonal, a classical algorithm for solving (1.9) is the projected gradient method (see Section 5.2), which performs a sequence of projections (1.18) for different values of β . Note that how to solve (1.18) efficiently is well studied; it can be achieved in $O(p)$ operations with a divide-and-conquer strategy [Brucker, 1984, Duchi et al., 2008].

In Figure 1.6, we illustrate the effect of the ℓ_1 -norm projection and compare it to the case of the ℓ_2 -norm. The corners of the ℓ_1 -ball are on the main axes and correspond to sparse solutions. Two of them are represented by red and green dots, with respective coordinates $(\mu, 0)$ and $(0, \mu)$. Most strikingly, a large part of the space in the figure, represented by red and green regions, ends up on these corners after projection. In contrast, the set of points that is projected onto the blue dot, is simply the blue line. The blue dot corresponds in fact to a dense solution with coordinates $(\mu/2, \mu/2)$. Therefore, the figure illustrates that the ℓ_1 -ball in two dimensions encourages solutions to be on its corners. In the case of the ℓ_2 -norm, the ball is isotropic, and treats every direction equally. In Figure 1.7, we represent these two balls in three dimensions, where we can make similar observations.

More formally, we can mathematically characterize our remarks about Figure 1.6. Consider a point \mathbf{y} in \mathbb{R}^p on the surface of the ℓ_1 -ball of radius $\mu = 1$, and define the set $\mathcal{N} \triangleq \{\mathbf{z} \in \mathbb{R}^p : \pi(\mathbf{z}) = \mathbf{y}\}$, where π is the projection operator onto the ℓ_1 -ball. Examples of pairs $(\mathbf{y}, \mathcal{N})$ have been presented in Figure 1.6; for instance, when \mathbf{y} is the red or green dot, \mathcal{N} is respectively the red or green region. It is particularly informative to study how \mathcal{N} varies with \mathbf{y} , which is the focus of the next proposition.

Proposition 1.1 (Characterization of the set \mathcal{N}).

For a non-zero vector \mathbf{y} in \mathbb{R}^p , the set \mathcal{N} defined in the previous paragraph can be written as $\mathcal{N} = \mathbf{y} + \mathcal{K}$, where \mathcal{K} is a polyhedral cone of dimension $p - \|\mathbf{y}\|_0 + 1$.

Proof. A classical theorem [see Bertsekas, 1999, Proposition B.11] allows us to rewrite \mathcal{N} as

$$\mathcal{N} = \{\mathbf{z} \in \mathbb{R}^p : \forall \|\mathbf{x}\|_1 \leq 1, (\mathbf{z} - \mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \leq 0\} = \mathbf{y} + \mathcal{K},$$

where $\mathbf{y} + \mathcal{K}$ denotes the Minkowski sum $\{\mathbf{y} + \mathbf{z} : \mathbf{z} \in \mathcal{K}\}$ between the set $\{\mathbf{y}\}$ and the cone \mathcal{K} defined as

$$\mathcal{K} \triangleq \{\mathbf{d} \in \mathbb{R}^p : \forall \|\mathbf{x}\|_1 \leq 1, \mathbf{d}^\top (\mathbf{x} - \mathbf{y}) \leq 0\}.$$

Note that in the optimization literature, \mathcal{K} is often called the “normal cone” to the unit ℓ_1 -ball at the point \mathbf{y} [Borwein and Lewis, 2006]. Equivalently, we have

$$\begin{aligned} \mathcal{K} &= \{\mathbf{d} \in \mathbb{R}^p : \max_{\|\mathbf{x}\|_1 \leq 1} \mathbf{d}^\top \mathbf{x} \leq \mathbf{d}^\top \mathbf{y}\} \\ &= \{\mathbf{d} \in \mathbb{R}^p : \|\mathbf{d}\|_\infty \leq \mathbf{d}^\top \mathbf{y}\}, \end{aligned} \quad (1.19)$$

where we have used the fact that quantity $\max_{\|\mathbf{x}\|_1 \leq 1} \mathbf{d}^\top \mathbf{x}$, called the dual-norm of the ℓ_1 -norm, is equal to $\|\mathbf{d}\|_\infty$ [see Bach et al., 2012a]. Note now that according to Hölder’s inequality, we also have $\mathbf{d}^\top \mathbf{y} \leq \|\mathbf{d}\|_\infty \|\mathbf{y}\|_1 \leq \|\mathbf{d}\|_\infty$ in Eq. (1.19). Therefore, the inequalities are in fact equalities. It is then easy to characterize vectors \mathbf{d} such that $\mathbf{d}^\top \mathbf{y} = \|\mathbf{d}\|_\infty \|\mathbf{y}\|_1$ and it is possible to show that \mathcal{K} is simply the set of vectors \mathbf{d} satisfying $\mathbf{d}[i] = \text{sign}(\mathbf{y}[i]) \|\mathbf{d}\|_\infty$ for all i such that $\mathbf{y}[i] \neq 0$.

This would be sufficient to conclude the proposition, but it is also possible to pursue the analysis and exactly characterize \mathcal{K} by finding a set of generators.⁴ Let us define the vector \mathbf{s} in $\{-1, 0, +1\}^p$ that carries the sparsity pattern of \mathbf{y} , more precisely, with $\mathbf{s}[i] = \text{sign}(\mathbf{y}[i])$ for all i such that $\mathbf{y}[i] \neq 0$, and $\mathbf{s}[i] = 0$ otherwise. Let us also define the set of indices $\{i_1, \dots, i_l\}$ corresponding to the l zero entries of \mathbf{y} , and \mathbf{e}_i in \mathbb{R}^p the binary vector whose entries are all zero but the i -th one that is equal to 1. Then, after a short calculation, we can geometrically characterize the polyhedral cone \mathcal{K} :

$$\mathcal{K} = \text{cone}(\mathbf{s}, \mathbf{s} - \mathbf{e}_{i_1}, \mathbf{s} + \mathbf{e}_{i_1}, \mathbf{s} - \mathbf{e}_{i_2}, \mathbf{s} + \mathbf{e}_{i_2}, \dots, \mathbf{s} - \mathbf{e}_{i_l}, \mathbf{s} + \mathbf{e}_{i_l}),$$

where the notation “cone” is defined in footnote 4. \square

It is now easy to see that the set \mathcal{K} “grows” with the number l of zero entries in \mathbf{y} , and that \mathcal{K} lives in a subspace of dimension $l + 1$ for

⁴A collection of vectors $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_l$ are called generators for a cone \mathcal{K} when \mathcal{K} consists of all positive combinations of the vectors \mathbf{z}_i . In other words, $\mathcal{K} = \{\sum_{i=1}^l \alpha_i \mathbf{z}_i : \alpha_i \geq 0\}$. In that case, we use the notation $\mathcal{K} = \text{cone}(\mathbf{z}_1, \dots, \mathbf{z}_l)$.

all non-zero vector \mathbf{y} . For example, when $l = 0$ —that is, \mathbf{y} is a dense vector (*e.g.*, the blue point in Figure 1.6(a)), \mathcal{K} is simply a half-line.

To conclude, the geometrical intuition to gain from this section is that *the Euclidean projection onto a convex set encourages solutions on singular points, such as edges or corners for polytopes*. Such a principle indeed applies beyond the ℓ_1 -norm. For instance, we illustrate the regularization effect of the ℓ_∞ -norm in Figure 1.8, whose corners coordinates have same magnitude.

Non-convex regularization. Even though it is well established that the ℓ_1 -norm encourages sparse solutions, it remains only a convex proxy of the ℓ_0 -penalty. Both in statistics and signal processing, other sparsity-inducing regularization functions have been proposed, in particular continuous relaxations of ℓ_0 that are non-convex [Frank and Friedman, 1993, Fan and Li, 2001, Daubechies et al., 2010, Gasso et al., 2009]. These functions are using a non-decreasing concave function $\varphi : \mathbb{R}^+ \mapsto \mathbb{R}$, and the sparsity-inducing penalty is defined as

$$\psi(\boldsymbol{\alpha}) \triangleq \sum_{i=1}^p \varphi(|\boldsymbol{\alpha}[i]|).$$

For example, the ℓ_q -penalty uses $\varphi : x \mapsto x^q$ [Frank and Friedman, 1993], or an approximation $\varphi : x \mapsto (x + \varepsilon)^q$; the reweighted- ℓ_1 algorithm of Fazel [2002], Fazel et al. [2003], Candès et al. [2008] implicitly uses $\varphi : x \mapsto \log(x + \varepsilon)$. These penalties typically lead to intractable estimation problems, but approximate solutions can be obtained with continuous optimization techniques (see Section 5.3).

The sparsity-inducing effect of the penalties ψ is known to be stronger than ℓ_1 . As shown in Figure 1.9(a), the magnitude of the derivative of φ grows when one approaches zero because of its concavity. Thus, in the one-dimensional case, ψ can be interpreted as *a force driving α towards the origin with increasing intensity when α gets closer to zero*. In terms of geometry, we also display the ℓ_q -ball in Figure 1.9(b), with the same red, blue, and green dots as in Figure 1.6. The part of the space that is projected onto the corners of the ℓ_q -ball is larger than that for ℓ_1 . Interestingly, the geometrical structure of the red and green regions are also more complex. Their combinatorial

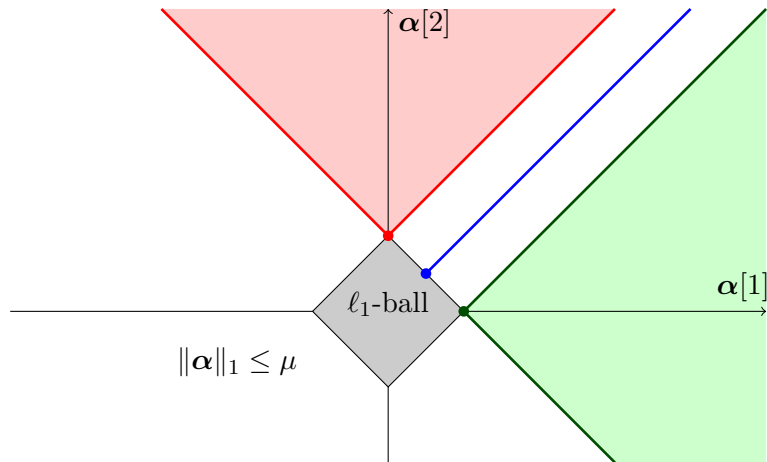
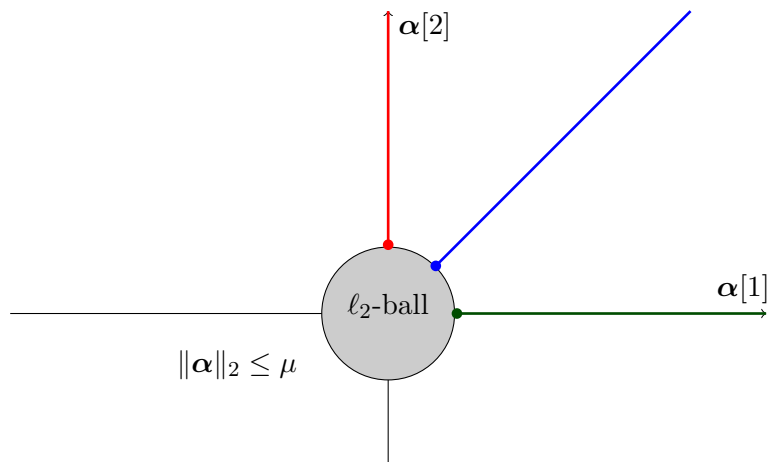
(a) Effect of the Euclidean projection onto the ℓ_1 -ball.(b) Effect of the Euclidean projection onto the ℓ_2 -ball.

Figure 1.6: Illustration in two dimensions of the projection operator onto the ℓ_1 -ball in Figure (a) and ℓ_2 -ball in Figure (b). The balls are represented in gray. All points from the red regions are projected onto the point of coordinates $(0, \mu)$ denoted by a red dot. Similarly, the green and blue regions are projected onto the green and blue dots, respectively. For the ℓ_1 -norm, a large part of the figure is filled by the red and green regions, whose points are projected to a sparse solution corresponding to a corner of the ball. For the ℓ_2 -norm, this is not the case: any non-sparse point—say, for instance on the blue line—is projected onto a non-sparse solution.

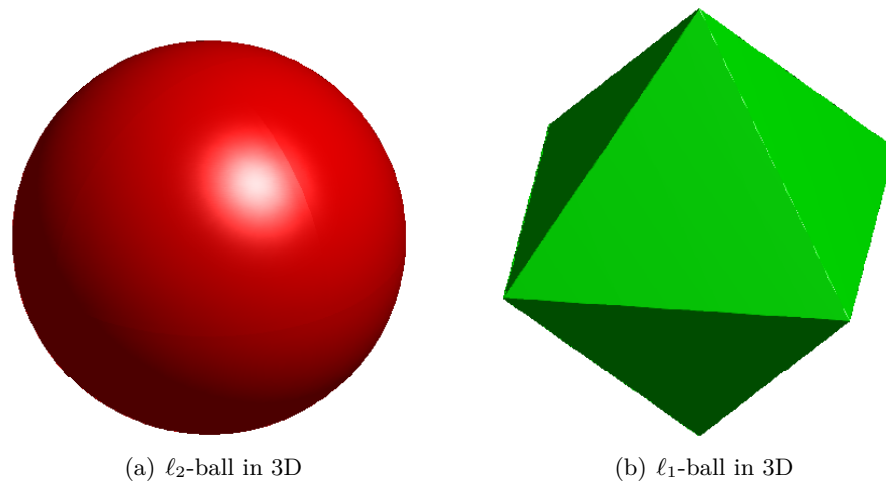


Figure 1.7: Representation in three dimensions of the ℓ_1 - and ℓ_2 -balls. Figure borrowed from Bach et al. [2012a], produced by Guillaume Obozinski.

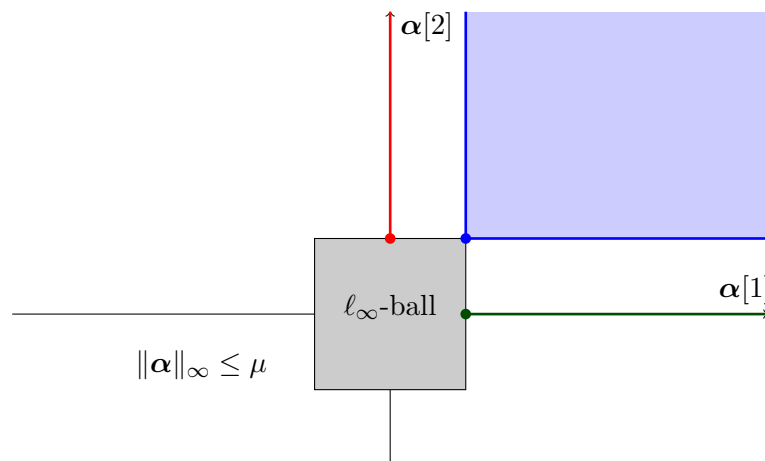


Figure 1.8: Similar illustration as Figure 1.6 for the ℓ_{∞} -norm. The regularization effect encourages solution to be on the corners of the ball, corresponding to points with the same magnitude $|\boldsymbol{\alpha}[1]| = |\boldsymbol{\alpha}[2]| = \mu$.

nature makes the projection problem onto the ℓ_q -ball more involved when $q < 1$.

The elastic-net. To cope with some instability issues of the estimators obtained with the ℓ_1 -regularization, Zou and Hastie [2005] have proposed to combine the ℓ_1 - and ℓ_2 -norms with a penalty called elastic-net:

$$\psi(\boldsymbol{\alpha}) \triangleq \|\boldsymbol{\alpha}\|_1 + \gamma\|\boldsymbol{\alpha}\|_2^2.$$

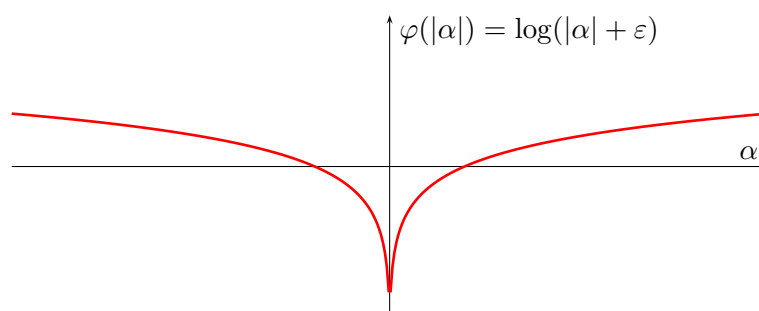
The effect of this penalty is illustrated in Figure 1.10. Compared to Figure 1.6, we observe that the red and green regions are smaller for the elastic-net penalty than for ℓ_1 . The sparsity-inducing effect is thus less aggressive than the one obtained with ℓ_1 .

Total variation. The anisotropic total variation penalty [Rudin et al., 1992] for one dimensional signals is simply the ℓ_1 -norm of finite differences

$$\psi(\boldsymbol{\alpha}) \triangleq \sum_{i=1}^{p-1} |\boldsymbol{\alpha}[i+1] - \boldsymbol{\alpha}[i]|,$$

which encourages piecewise constant signals. It is also known in statistics under the name of “fused Lasso” [Tibshirani et al., 2005]. The penalty can easily be extended to two-dimensional signals, and has been widely used for regularizing inverse problems in image processing [Chambolle, 2005].

Group sparsity. In some cases, variables are organized into predefined groups forming a partition \mathcal{G} of $\{1, \dots, p\}$, and one is looking for a solution $\boldsymbol{\alpha}^*$ such that *variables belonging to the same group of \mathcal{G} are set to zero together*. For example, such groups have appeared in Section 1.2 about wavelets, where \mathcal{G} could be defined according to neighborhood relationships of wavelet coefficients. Then, when it is known beforehand that a problem solution only requires a few groups of variables to explain the data, a regularization function automatically selecting the relevant groups has been shown to improve the prediction performance or the interpretability of the solution [Turlach et al., 2005,



(a) Illustration of a non-convex sparsity-inducing penalty.

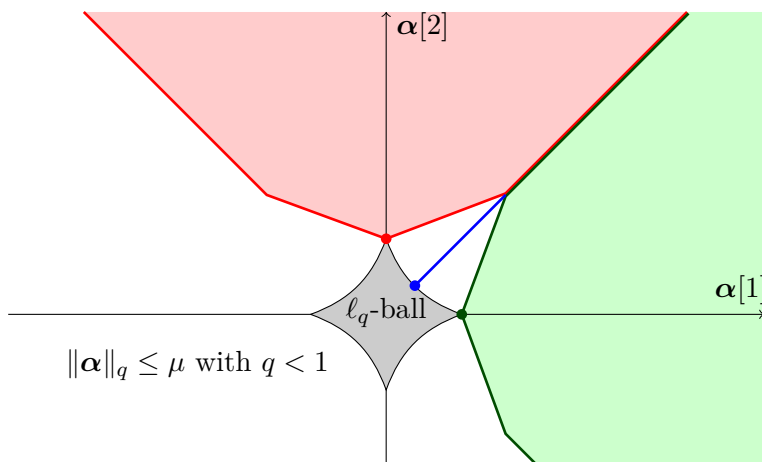
(b) ℓ_q -ball with $q < 1$.

Figure 1.9: Illustration of the sparsity-inducing effect of a non-convex penalty. In (a), we plot the non-convex penalty $\alpha \mapsto \log(|\alpha| + \varepsilon)$, and in (b), we present a similar figure as 1.6 for the ℓ_q -penalty, when choosing $q < 1$.

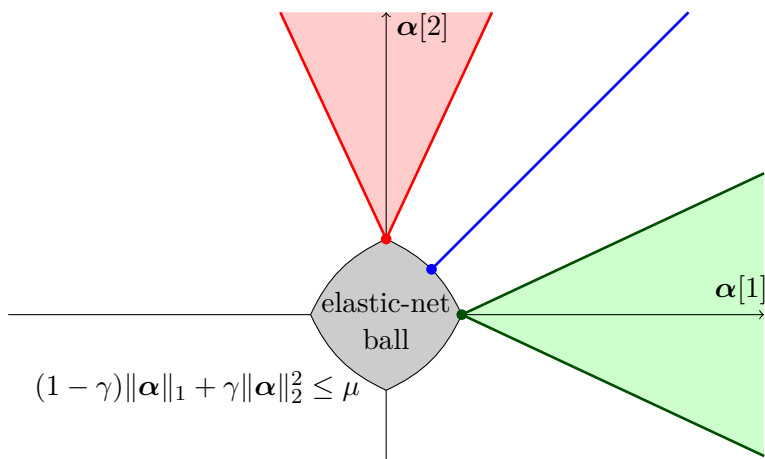


Figure 1.10: Similar figure as 1.6 for the elastic-net penalty.

Yuan and Lin, 2006, Obozinski et al., 2009, Huang and Zhang, 2010]. The group sparsity principle is illustrated in Figure 1.11(b).

An appropriate regularization function to obtain a group-sparsity effect is known as “Group-Lasso” penalty and is defined as

$$\psi(\boldsymbol{\alpha}) = \sum_{g \in \mathcal{G}} \|\boldsymbol{\alpha}[g]\|_q, \quad (1.20)$$

where $\|\cdot\|_q$ is either the ℓ_2 or ℓ_∞ -norm. To the best of our knowledge, such a penalty appears in the early work of Grandvalet and Canu [1999] and Bakin [1999] for $q = 2$, and Turlach et al. [2005] for $q = \infty$. It has been popularized later by Yuan and Lin [2006].

The function ψ in (1.20) is a norm, thus convex, and can be interpreted as the ℓ_1 -norm of the vector $[\|\boldsymbol{\alpha}[g]\|_q]_{g \in \mathcal{G}}$ of size $|\mathcal{G}|$. Consequently, the sparsity-inducing effect of the ℓ_1 -norm is applied at the group level. The penalty is highly related to the group-thresholding approach for wavelets, since the group-thresholding estimator (1.7) is linked to ψ through Eq. (1.8).

In Figure 1.13(a), we visualize the unit ball of a Group-Lasso norm obtained when \mathcal{G} contains two groups $\mathcal{G} = \{\{1, 2\}, \{3\}\}$. The ball has two singularities: the top and bottom corners, corresponding to solutions where variables 1 and 2 are simultaneously set to zero, and the

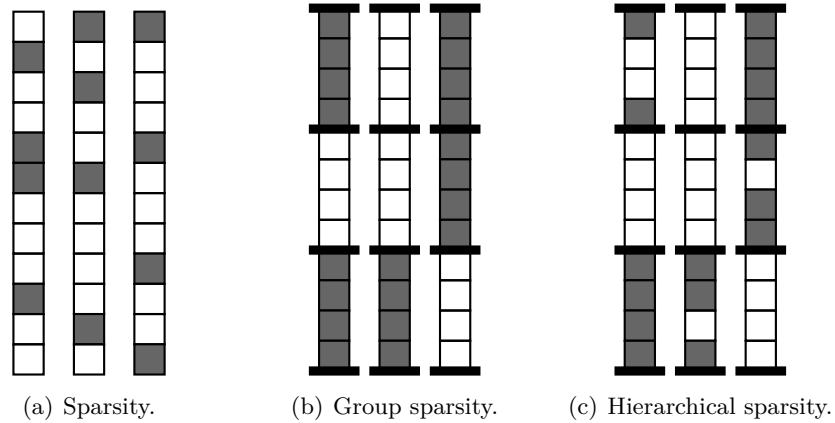


Figure 1.11: Illustration of the sparsity, group sparsity, and hierarchical sparsity principles. Each column represents the sparsity pattern of a vector with 12 variables. Non-zero coefficients are represented by gray squares. On the left (a), the vectors are obtained with a simple sparsity-inducing penalty, such as the ℓ_1 -norm, and the non-zero variables are scattered. In the middle figure (b), a group sparsity-inducing penalty with three groups of variables is used. The sparsity patterns respect the group structure. On the right (c), we present the results obtained with a hierarchical penalty consisting of the Group Lasso plus the ℓ_1 -norm. The group structure is globally respected, but some variables within a group can be discarded.

middle circle, corresponding to solutions where variable 3 only is set to zero. As expected, the geometry of the ball induces the group-sparsity effect.

Structured sparsity. Group-sparsity is a first step towards the more general idea that a regularization function can encourage sparse solutions with a particular structure. This notion is called *structured sparsity* and has been introduced under a large number of different points of views [Zhao et al., 2009, Jacob et al., 2009, Jenatton et al., 2011a, Baraniuk et al., 2010, Huang et al., 2011]. To some extent, it follows the concept of group-thresholding introduced in the wavelet literature, which we have presented in Section 1.2. In this paragraph, we briefly review some of these works, but for a more detailed review, we refer the reader to [Bach et al., 2012b].

Some penalties are non-convex. For instance, Huang et al. [2011] and Baraniuk et al. [2010] propose two different combinatorial approaches based on a predefined set \mathcal{G} of possibly overlapping groups of variables. These penalties encourage solutions whose support is in the *union of a few number groups*, but they lead to NP-hard optimization problems. Other penalties are convex. In particular, Jacob et al. [2009] introduce a sparsity-inducing norm that is exactly a convex relaxation of the penalty of Huang et al. [2011], even though these two approaches were independently developed at the same time. As a result, the convex penalty of Jacob et al. [2009] encourages a similar structure as the one of Huang et al. [2011].

By following a different direction, the Group-Lasso penalty (1.20) has been considered when the groups are allowed to overlap [Zhao et al., 2009, Jenatton et al., 2011a]. As a consequence, variables belonging to the same groups are encouraged to be set to zero together. It was proposed for hierarchical structures by Zhao et al. [2009] with the following rule: whenever two groups g and h are in \mathcal{G} , they should be either disjoint, or one should be included in another. Examples of such hierarchical group structures are given in Figures 1.11 and 1.12. The effect of the penalty is to encourage sparsity patterns that are rooted subtrees. Equivalently, a variable can be non-zero only

if its parent in the tree is non-zero, which is the main property of the zero-tree coding scheme introduced in the wavelet literature [Shapiro, 1993], and already illustrated in Figure 1.2.

Finally, Jenatton et al. [2011a] has extended the hierarchical penalty of Zhao et al. [2009] to more general group structures, for example when variables are organized on a two-dimensional grid, encouraging neighbor variables to be simultaneously set to zero. We conclude this brief presentation of structured sparsity with Figure 1.13, where we present the unit balls of some sparsity-inducing norms. Each of them exhibits singularities and encourages particular sparsity patterns.

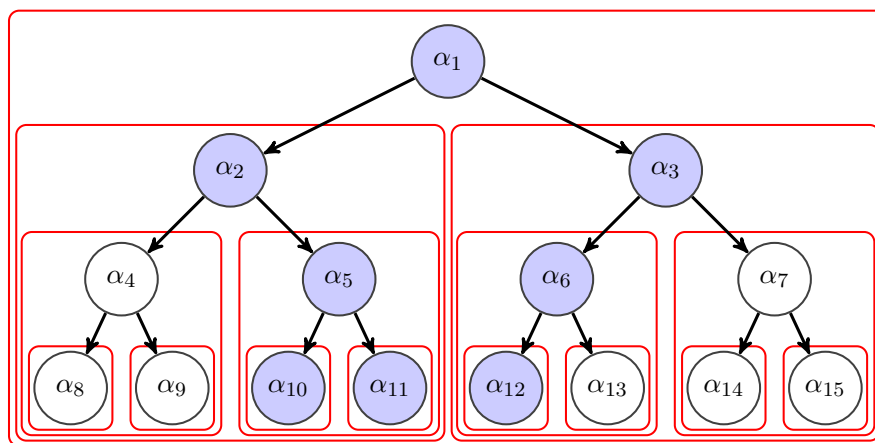


Figure 1.12: Illustration of the hierarchical sparsity of Zhao et al. [2009], which generalizes the zero-tree coding scheme of Shapiro [1993]. The groups of variables correspond to the red rectangles. The empty nodes represent variables that are set to zero. They are contained in three groups: $\{4, 8, 9\}$, $\{13\}$, $\{7, 14, 15\}$.

Spectral sparsity. Another form of parsimony has been devised in the spectral domain [Fazel et al., 2001, Srebro et al., 2005]. For estimation problems where model parameters are matrices, the rank has been used as a natural regularization function. The rank of a matrix is equal to the number of non-zero singular values, and thus, it can be interpreted as the ℓ_0 -penalty of the matrix spectrum. Unfortunately, due to the combinatorial nature of ℓ_0 , the rank penalization typically leads to

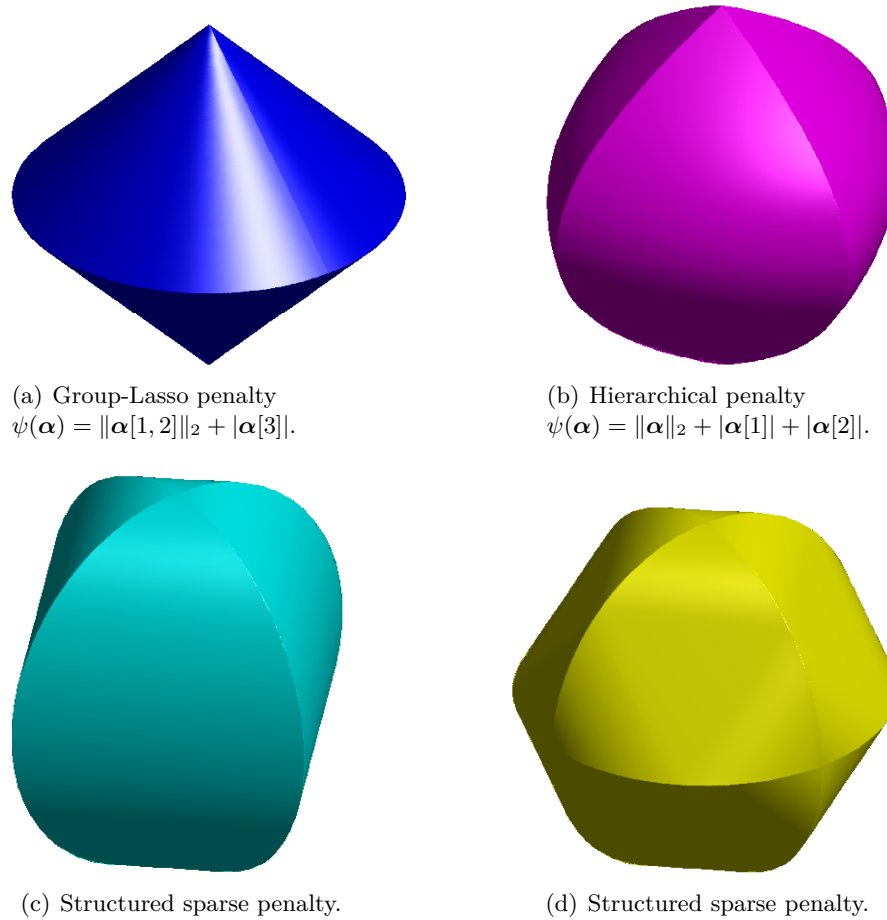


Figure 1.13: Visualization in three dimensions of unit balls corresponding to various sparsity-inducing norms. (a): Group Lasso penalty; (b): hierarchical penalty of Zhao et al. [2009]; (c) and (d): examples of structured sparsity-inducing penalties of Jacob et al. [2009]. Figure borrowed from Bach et al. [2012a], produced by Guillaume Obozinski.

intractable optimization problems.

A natural convex relaxation has been introduced in the control theory literature by Fazel et al. [2001] and consists of computing the ℓ_1 -norm of the spectrum—that is, simply the sum of the singular values. The resulting penalty appears under different names, the most common ones being the trace, nuclear, or Schatten norm. It is defined for a matrix \mathbf{A} in $\mathbb{R}^{p \times k}$ with $k \geq p$ as

$$\|\mathbf{A}\|_* \triangleq \sum_{i=1}^p s_i(\mathbf{A}),$$

where $s_i(\mathbf{A})$ is the i -th singular value of \mathbf{A} . Traditional applications of the trace norm in machine learning are matrix completion or collaborative filtering [Pontil et al., 2007, Abernethy et al., 2009]. These problems have become popular with the need of scalable recommender systems for video streaming providers. The goal is to infer movie preferences for each customer, based on their partial movie ratings. Typically, the matrix is of size $p \times k$, where p is the number of movies and k is the number of users. Each user gives a score for a few movies, corresponding to some entries of the matrix, and the recommender system tries to infer the missing values. Similar techniques have also recently been used in other fields, such as in genomics to infer missing genetic information [Chi et al., 2013].

1.4 Dictionary learning

We have previously presented various formulations where a signal \mathbf{x} in \mathbb{R}^m is approximated by a sparse linear combination of a few columns of a matrix \mathbf{D} in $\mathbb{R}^{m \times p}$. In the context of signal and image processing, this matrix is often called *dictionary* and its columns *atoms*. As seen in Section 1.2, a large amount of work has been devoted in the wavelet literature for designing a good dictionary adapted to natural images.

In neuroscience, Olshausen and Field [1996, 1997] have proposed a significantly different approach to sparse modeling consisting of adapting the dictionary to training data. Because the size of natural images is too large for learning a full matrix \mathbf{D} , they have chosen to learn the dictionary on natural image patches, *e.g.*, of size $m = 16 \times 16$

pixels, and have demonstrated that their method could automatically discover interpretable structures. We discuss this topic in more details in Section 2.

The motivation of Olshausen and Field [1996, 1997] was to show that the structure of natural images is related to classical theories of the mammalian visual cortex. Later, dictionary learning found numerous applications in image restoration, and was shown to significantly outperform off-the-shelf bases for signal reconstruction [see, *e.g.*, Elad and Aharon, 2006, Mairal et al., 2008c, 2009, Protter and Elad, 2009, Yang et al., 2010a].

Concretely, given a dataset of n training signals $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, dictionary learning can be formulated as the following minimization problem

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda\psi(\boldsymbol{\alpha}_i), \quad (1.21)$$

where $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n]$ carries the decomposition coefficients of the signals $\mathbf{x}_1, \dots, \mathbf{x}_n$, ψ is sparsity-inducing regularization function, and \mathcal{C} is typically chosen as the following set:

$$\mathcal{C} \triangleq \{\mathbf{D} \in \mathbb{R}^{m \times p} : \forall j \quad \|\mathbf{d}_j\|_2 \leq 1\}.$$

To be more precise, Olshausen and Field [1996] proposed several choices for ψ ; their experiments were for instance conducted with the ℓ_1 -norm, or with the smooth function $\psi(\boldsymbol{\alpha}) \triangleq \sum_{j=1}^p \log(\varepsilon + \boldsymbol{\alpha}[j]^2)$, which has an approximate sparsity-inducing effect. The constraint $\mathbf{D} \in \mathcal{C}$ was also not explicitly modeled in the original dictionary learning formulation; instead, the algorithm of Olshausen and Field [1996] includes a mechanism to control and rescale the ℓ_2 -norm of the dictionary elements. Indeed, without such a mechanism, the norm of \mathbf{D} would arbitrarily go to infinity, leading to small values for the coefficients $\boldsymbol{\alpha}_i$ and making the penalty ψ ineffective.

The number of samples n is typically large, whereas the signal dimension m is small. The number of dictionary elements p is often chosen larger than m —in that case, the dictionary is said to be *over-complete*—even though a choice $p < m$ often leads to reasonable results in many applications. For instance, a typical setting would be to

have $m = 10 \times 10$ pixels for natural image patches, a dictionary of size $p = 256$, and more than 100 000 training patches.

A large part of this monograph is related to dictionary learning and thus we only briefly discuss this matter in this introduction. Section 2 is indeed devoted to unsupervised learning techniques for natural image patches, including dictionary learning; Sections 3 and 4 present a large number of applications in image processing and computer vision; how to solve (1.21) is explained in Section 5.5 about optimization.

Matrix factorization point of view. An equivalent representation of (1.21) is the following regularized matrix factorization problem

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{D}\mathbf{A}\|_{\text{F}}^2 + \lambda \Psi(\mathbf{A}), \quad (1.22)$$

where $\Psi(\mathbf{A}) = \sum_{i=1}^n \psi(\boldsymbol{\alpha}_i)$. Even though reformulating (1.21) as (1.22) is simply a matter of using different notation, seeing dictionary learning as a matrix factorization problem opens up interesting perspectives. In particular, it makes obvious some links with other unsupervised learning approaches such as non-negative matrix factorization [Paatero and Tapper, 1994], clustering techniques such as K-means, and others [see Mairal et al., 2010a]. These links will be further developed in Section 2.

Risk minimization point of view. Dictionary learning can also be seen from a machine learning point of view. Indeed, dictionary learning can be written as

$$\min_{\mathbf{D} \in \mathcal{C}} \left\{ f_n(\mathbf{D}) \triangleq \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{D}) \right\},$$

where $L : \mathbb{R}^m \times \mathbb{R}^{m \times p}$ is a loss function defined as

$$L(\mathbf{x}, \mathbf{D}) \triangleq \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \psi(\boldsymbol{\alpha}).$$

The quantity $L(\mathbf{x}, \mathbf{D})$ should be small if \mathbf{D} is “good” at representing the signal \mathbf{x} in a sparse fashion, and large otherwise. Then, $f_n(\mathbf{D})$ is called the *empirical cost*.

However, as pointed out by Bottou and Bousquet [2008], one is usually not interested in the exact minimization of the empirical cost $f_n(\mathbf{D})$ for a fixed n , which may lead to overfitting on the training data, but instead in the minimization of the *expected cost*, which measures the quality of the dictionary on new unseen data:

$$f(\mathbf{D}) \triangleq \mathbb{E}_{\mathbf{x}}[L(\mathbf{x}, \mathbf{D})] = \lim_{n \rightarrow \infty} f_n(\mathbf{D}) \text{ a.s.},$$

where the expectation is taken relative to the (unknown) probability distribution of the data.⁵

The expected risk minimization formulation is interesting since it paves the way to stochastic optimization techniques when a large amount of data is available [Mairal et al., 2010a] and to theoretical analysis [Maurer and Pontil, 2010, Vainsencher et al., 2011, Gribonval et al., 2013], which are developed in Sections 5.5 and 1.6, respectively.

Constrained variants. Following the original formulation of Olshausen and Field [1996, 1997], we have chosen to present dictionary learning where the regularization function is used as a penalty, even though it can also be used as a constraint as in (1.9). Then, natural variants of (1.21) are

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \text{ s.t. } \psi(\boldsymbol{\alpha}_i) \leq \mu. \quad (1.23)$$

or

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \sum_{i=1}^n \psi(\boldsymbol{\alpha}_i) \text{ s.t. } \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \leq \varepsilon. \quad (1.24)$$

Note that (1.23) and (1.24) are not equivalent to (1.21). For instance, problem (1.23) can be reformulated using a Lagrangian function [Boyd and Vandenberghe, 2004] as

$$\min_{\mathbf{D} \in \mathcal{C}} \sum_{i=1}^n \left(\max_{\lambda_i \geq 0} \min_{\boldsymbol{\alpha}_i \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda_i (\psi(\boldsymbol{\alpha}_i) - \mu) \right),$$

where the optimal λ_i 's are not necessarily equal to each other, and their relation with the constraint parameter μ is unknown in advance.

⁵We use ‘‘a.s.’’ to denote almost sure convergence.

A similar discussion can be conducted for (1.24) and it is thus important in practice to choose one of the formulations (1.21), (1.23), or (1.24); the best one depends on the problem at hand and there is no general rule for preferring one instead of another.

1.5 Compressed sensing and sparse recovery

Finally, we conclude our historical tour of parsimony with recent theoretical results obtained in signal processing and statistics. We focus on methods based on the ℓ_1 -norm, *i.e.*, the basis pursuit formulation of (1.9)—more results on structured sparsity-inducing norms are presented by Bach et al. [2012b].

Most analyses rely on particular assumptions regarding the problem. We start this section with a cautionary note from Hocking [1976]:

The problem of selecting a subset of independent or predictor variables is usually described in an idealized setting. That is, it is assumed that (a) the analyst has data on a large number of potential variables which include all relevant variables and appropriate functions of them plus, possibly, some other extraneous variables and variable functions and (b) the analyst has available “good” data on which to base the eventual conclusions. In practice, the lack of satisfaction of these assumptions may make a detailed subset selection analysis a meaningless exercise.

In this section, we present such theoretical results where the assumptions are often not met in practice, but also results that either (1) can have an impact on the practice of sparse recovery or (2) do not need strong assumptions.

From support recovery to signal denoising. Given a signal \mathbf{x} in \mathbb{R}^m and a dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$ with ℓ_2 -normalized columns, throughout this section, we assume that \mathbf{x} is generated as $\mathbf{x} = \mathbf{D}\boldsymbol{\alpha}^* + \boldsymbol{\varepsilon}$ with a sparse vector $\boldsymbol{\alpha}^*$ in \mathbb{R}^p and an additive noise $\boldsymbol{\varepsilon}$ in \mathbb{R}^m . For simplicity, we consider $\boldsymbol{\alpha}^*$ and \mathbf{D} as being deterministic while the noise is

random, independent and identically distributed, with zero mean and finite variance σ^2 .

The different formulations presented earlier in Section 1.3, for instance basis pursuit, provide estimators $\hat{\boldsymbol{\alpha}}$ of the “true” vector $\boldsymbol{\alpha}^*$. Then, the three following goals have been studied in sparse recovery, typically in decreasing order of hardness:

- **support recovery and sign consistency:** we want the support of $\hat{\boldsymbol{\alpha}}$ (*i.e.*, the set of non-zero elements) to be the same or to be close to the one of $\boldsymbol{\alpha}^*$. The problem is often called “model selection” in statistics and “support recovery” in signal processing; it is often refined to the estimation of the full sign pattern—that is, among the non-zero elements, we also want the correct sign to be estimated.
- **code estimation:** the distance $\|\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*\|_2$ should be small. In statistical terms, this correspond to the “estimation” of $\boldsymbol{\alpha}^*$.
- **signal denoising:** regardless of code estimation, we simply want the distance $\|\mathbf{D}\hat{\boldsymbol{\alpha}} - \mathbf{D}\boldsymbol{\alpha}^*\|_2$ to be small; the goal is not to obtain exactly $\boldsymbol{\alpha}^*$, but simply to obtain a good denoised version of the signal $\mathbf{x} = \mathbf{D}\boldsymbol{\alpha}^* + \boldsymbol{\varepsilon}$.

A good code estimation performance does imply a good denoising performance but the converse is not true in general. In most analyses, support recovery is harder than code estimation. As detailed below, the sufficient conditions for good support recovery lead indeed to good estimation.

High-dimensional phenomenon. Without sparsity assumptions, even the simplest denoising task can only achieve denoising errors of the order $\frac{1}{n}\|\mathbf{D}\hat{\boldsymbol{\alpha}} - \mathbf{D}\boldsymbol{\alpha}^*\|_2^2 \approx \frac{\sigma^2 p}{m}$, which is attained for ordinary least-squares, and is the best possible [Tsybakov, 2003]. Thus, in order to have at least a good denoising performance (prediction performance in statistics), either the noise σ is small, or the signal dimension m (the number of samples) is much larger than the number of atoms p (the number of variables to select from).

When making the assumption that the true code α^* is sparse with at most k non zeros, smaller denoising errors can be obtained. In that case, it is possible indeed to replace the scaling $\frac{\sigma^2 p}{m}$ by $\frac{\sigma^2 k \log p}{m}$. Thus, even when p is much larger than m , as long as $\log p$ is much smaller than m , we may have good prediction performance. However, this high-dimensional phenomenon currently⁶ comes at a price: (1) either an exhaustive search over the subsets of size k needs to be performed [Massart, 2003, Bunea et al., 2007, Raskutti et al., 2011] or (2) some assumptions have to be made regarding the dictionary \mathbf{D} , which we now describe.

Sufficient conditions for high-dimensional fast rates. Most sufficient conditions have the same flavor. A dictionary behaves well if the off-diagonal elements of $\mathbf{D}^\top \mathbf{D}$ are small, in other words, if there is little correlation between atoms. However, the notion of coherence (the maximal possible correlation between two atoms) was the first to emerge [see, e.g., Elad and Bruckstein, 2002, Gribonval and Nielsen, 2003], but it is not sufficient to obtain a high-dimensional phenomenon.

In the noiseless setting, Candes and Tao [2005] and Candès et al. [2006] introduced the *restricted isometry property* (RIP), which states that all submatrices of size $k \times k$ of $\mathbf{D}^\top \mathbf{D}$ should be close to isometries, that is, should have all of their eigenvalues sufficiently close to one. With such an assumption, the Lasso behaves well: it recovers the true support and estimates the code α^* and the signal $\mathbf{D}\alpha^*$ with an error of order $\frac{\sigma^2 k \log p}{m}$.

The main advantage of the RIP assumption is that one may exhibit dictionaries for which it is satisfied, usually obtained by normalizing a matrix \mathbf{D} obtained from independent Gaussian entries, which may satisfy the condition that $(k \log p)/m$ remains small. Thus, the sufficient conditions are not vacuous. However, the RIP assumption has two main drawbacks: first, it cannot be checked on a given dictionary \mathbf{D} without checking all $O(p^k)$ submatrices of size k ; second, it may be weakened if the goal is support recovery or simply estimation performance (code

⁶Note that recent research suggests that this fast rate of $\frac{\sigma^2 k \log p}{m}$ cannot be achieved by polynomial-time algorithms [Zhang et al., 2014].

recovery).

There is therefore a need for sufficient conditions that can be checked in polynomial time while ensuring sparse recovery. However, none currently exists with the same scalings between k , p and m [see, *e.g.*, Juditsky and Nemirovski, 2011, d’Aspremont and El Ghaoui, 2011]. When refining to support recovery, Fuchs [2005], Tropp [2004], Wainwright [2009] provide sufficient and necessary conditions of a similar flavor than requiring that all submatrices of size k are sufficiently close to orthogonal. For the tightest conditions, see, *e.g.*, Bühlmann and Van De Geer [2011]. Note that these conditions are also typically sufficient for algorithms that are not based explicitly on convex optimization [Tropp, 2004].

Finally, it is important to note that (a) most of the theoretical results advocate a value for the regularization parameter λ proportional to $\sigma\sqrt{m\log p}$, which unfortunately depends on the noise level σ (which is typically unknown in practice), and that (b) for orthogonal dictionaries, all of these assumptions are met; however, this imposes $p = m$.

Compressed sensing vs. statistics. Our earlier quote from Hocking [1976] applies to sparse estimation as used in statistics for least-squares regression, where the dictionary \mathbf{D} is simply the input data and \mathbf{x} the output data. In most situations, there are some variables, represented by columns of \mathbf{D} , that are heavily correlated. Therefore, in most practical situations, the assumptions do not apply. However, it does not mean that the high-dimensional phenomenon does not apply in a weaker sense (see the next paragraph for slow rates); moreover it is important to remark that there are other scenarios, beyond statistical variable selection, where the dictionary \mathbf{D} may be chosen.

In particular, in signal processing, the dictionary \mathbf{D} may be seen as *measurements*—that is, we want to encode $\boldsymbol{\alpha}^*$ in \mathbb{R}^p using m linear measurements $\mathbf{D}\boldsymbol{\alpha}^*$ in \mathbb{R}^m for m much larger than p . What the result of Candes and Tao [2005] alluded to earlier shows is that for random measurements, one can recover a k -sparse $\boldsymbol{\alpha}^*$ from (a potentially noisy version of) $\mathbf{D}\boldsymbol{\alpha}^*$, with overwhelming probability, as long as $(k \log p)/m$ remains small. This is the core idea behind compressive sensing. See

more details from Donoho [2006], Candès and Wakin [2008].

High-dimensional slow rates. While sufficient conditions presented earlier are often not met beyond random dictionaries, for the basis pursuit/Lasso formulation from (1.9), the high-dimensional phenomenon may still be observed, but only for the denoising situation and with a weaker result. Namely, as shown by Greenshtein [2006] and Bühlmann and Van De Geer [2011, Corollary 6.1], *without assumptions regarding correlations*, we have $\frac{1}{n}\|\mathbf{D}\hat{\boldsymbol{\alpha}}-\mathbf{D}\boldsymbol{\alpha}^*\|_2^2 \approx \sqrt{\frac{\sigma^2 k^2 \log p}{m}}$. Note that this slower rate does not readily extend to non-convex formulations.

Impact on dictionary learning. The dictionary learning framework which we describe in this monograph relies on sparse estimation, that is, given the dictionary \mathbf{D} , the estimation of the code $\boldsymbol{\alpha}$ may be analyzed using the tools we have presented in this section. However, the dictionaries that are learned do not exhibit low correlations between atoms and thus theoretical results do not apply (see dedicated results in the next section). However, they suggest that (a) the codes $\boldsymbol{\alpha}$ may not be unique in general and caution has to be observed when representing a signal \mathbf{x} by its code $\boldsymbol{\alpha}$, (b) methods based on ℓ_1 -penalization are more robust as they still provably perform denoising in presence of strong correlations and (c) incoherence promoting may be used in order to obtain better-behaved dictionaries [see, e.g., Ramirez et al., 2009].

1.6 Theoretical results about dictionary learning

Dictionary learning, as formulated in Eq. (1.21), may be seen from several perspectives, mainly as an unsupervised learning or a matrix factorization problem. While the supervised learning problem from the previous section (sparse estimation of a single signal given the dictionary) comes with many theoretical analyses, there are still few theoretical results of the same kind for dictionary learning. In this section, we present some of them. For simplicity, we assume that we penalize

with the ℓ_1 -norm and consider the minimization of

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1, \quad (1.25)$$

where $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n]$ carries the decomposition coefficients of the signals $\mathbf{x}_1, \dots, \mathbf{x}_n$, and \mathcal{C} is chosen as the following set:

$$\mathcal{C} \triangleq \{\mathbf{D} \in \mathbb{R}^{m \times p} : \forall j \quad \|\mathbf{d}_j\|_2 \leq 1\}.$$

Non-convex optimization problem. After imposing parsimony through the ℓ_1 -norm, given \mathbf{D} the objective function is convex in $\boldsymbol{\alpha}$, given $\boldsymbol{\alpha}$ the objective and constraints are convex in \mathbf{D} . However, the objective function is not jointly convex, which is typical of unsupervised learning formulations. Hence, we consider an optimization problem for which it is not possible in general to guarantee that we are going to obtain the global minimum; the same applies to EM-based approaches [Dempster et al., 1977] or K-means [see, *e.g.*, Bishop, 2006].

Symmetries. Worse, the problem in Equation (1.25) exhibits several symmetries and admits multiple global optima, and the descent methods that are described in Section 5 will also have the same invariance property. For example, the columns of \mathbf{D} and rows of \mathbf{A} can be submitted to $p!$ arbitrary (but consistent) permutations. There are also sign ambiguities: in fact, if (\mathbf{D}, \mathbf{A}) is solution of (1.25), so is $(\mathbf{D}\text{diag}(\boldsymbol{\varepsilon}), \text{diag}(\boldsymbol{\varepsilon})\mathbf{A})$, where $\boldsymbol{\varepsilon}$ is a vector in $\{-1, +1\}^p$ that carries a sign pattern. Therefore, for every one of the $p!$ possible atom orders, the dictionary learning problem admits 2^p equivalent solutions. In other words, for a solution (\mathbf{D}, \mathbf{A}) , the pair $(\mathbf{D}\boldsymbol{\Gamma}, \boldsymbol{\Gamma}^{-1}\mathbf{A})$ is also solution, where $\boldsymbol{\Gamma}$ is a *generalized permutation* formed by the product of a diagonal matrix with $+1$ and -1 's on its diagonal with a permutation matrix (in particular, $\boldsymbol{\Gamma}$ is thus orthogonal).

The fact that there are no other transformations $\boldsymbol{\Gamma}$ such that $(\mathbf{D}\boldsymbol{\Gamma}, \boldsymbol{\Gamma}^{-1}\mathbf{A})$ is also solution of Eq. (1.25) for *all* solutions (\mathbf{D}, \mathbf{A}) of this problem follows from a general property of isometries of the ℓ_q norm for finite values of q such that $q \geq 1$ and $q \neq 2$ [Li and So, 1994].

A manifold interpretation of sparse coding with projective geometry.

The interpretation of sparse coding as a *locally* linear representation of a non-linear “manifold” is problematic because certain signals/features are best thought of as “points” in some space rather than vectors. For example, what does it mean to “add” two natural image patches? The simplest point structure that one can think of is affine or projective, and we show below that sparse coding indeed admits a natural interpretation in this setting, at least for normalized signals.

Indeed, Let us restrict our attention from now on to unit-norm signals, as is customary in image processing after the usual centering and normalization steps, which will be studied in Section 2.1.⁷ Note that the dictionary elements \mathbf{d}_j in a solution $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ of Eq. (1.25) also have unit norm by construction.

Let us now consider the “half sphere”

$$\mathbb{S}_+^{m-1} \triangleq \left\{ \mathbf{d} \in \mathbb{S}^{m-1} : \text{the first non-zero coefficient of } \mathbf{d} \text{ is positive} \right\}, \quad (1.26)$$

where \mathbb{S}^{m-1} is the unit sphere of dimension $m - 1$ formed by the unit vectors of \mathbb{R}^m .⁸ A direct consequence of the sign ambiguities of dictionary learning discussed in the previous paragraph is that, for any solution (\mathbf{D}, \mathbf{A}) of Eq. (1.25), there is an equivalent solution $(\mathbf{D}', \mathbf{A}')$ with all columns of \mathbf{D}' in \mathbb{S}_+^{m-1} . Indeed, suppose some column \mathbf{d}_j is not in \mathbb{S}_+^{m-1} , and let $\mathbf{d}_j[i]$ be its first non-zero coefficient (which is necessarily negative since $\mathbf{d}_j \notin \mathbb{S}_+^{m-1}$). We can replace \mathbf{d}_j by $-\mathbf{d}_j$ and the corresponding row of the matrix \mathbf{A} by its opposite to construct an equivalent minimum of the dictionary learning problem in $\mathbb{S}_+^{m-1} \times \mathbb{R}^{p \times n}$.

Likewise, we can restrict the signals \mathbf{x}_i to lie in \mathbb{S}_+^{m-1} since replacing \mathbf{x}_i by its opposite for a given dictionary simply amounts to replacing the code α_i by its opposite. Note that this identifies a patch with its “negative”, but remember that the sign of its code elements is not uniquely defined in the first place in conventional dictionary learning settings (it is uniquely defined if we insist that the dictionary elements

⁷Note that the fact that the individual signals are centered does not imply that the dictionary elements are.

⁸Similarly, one may define the set \mathbb{S}_-^{m-1} by replacing “positive” by “non-negative” in (1.26). The two sets \mathbb{S}_+^{m-1} and \mathbb{S}_-^{m-1} form a partition of \mathbb{S}^{m-1} with equal volume, and indeed, each one geometrically corresponds to a half sphere.

belong to \mathbb{S}_+^{m-1}). This allows us to identify both the dictionary elements and the signals with points in the projective space $\mathbb{P}^{m-1} = P(\mathbb{R}^m)$. Any k independent column vectors of \mathbf{D} define a $(k-1)$ -dimensional projective subspace of \mathbb{P}^{m-1} (see Figure 1.14).

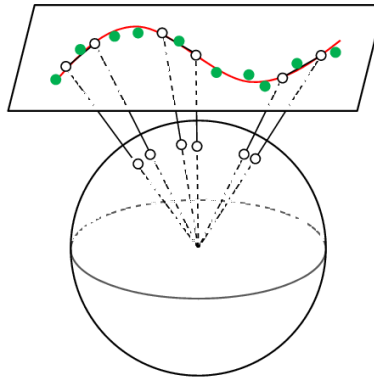


Figure 1.14: An illustration of the projective interpretation of sparse coding.

In particular, if the data signals are assumed to be sampled from a “noisy manifold” of dimension $k-1$ embedded in \mathbb{P}^{m-1} , an approximation of some sample \mathbf{x} by a sparse linear combination of k elements of \mathbf{D} can be thought of as lying in (or near) the $k-1$ dimensional “tangent plane” there.

Consistency results. Given the dictionary learning problem from a finite number of signals, there are several interesting theoretical questions to be answered. The first natural question is to understand the properties of the cost function that is minimized when the number of signals tends to infinity, and in particular how it converges to the expectation under the signal generating distribution [Vainsencher et al., 2011, Maurer and Pontil, 2010]. Then, given the non-convexity of the optimization problems, local consistency results may be obtained, by showing that the cost function which is minimized has a local mini-

mum around the pairs $(\mathbf{D}^*, \mathbf{A}^*)$ that has generated the data. Given RIP-based assumptions on the dictionary \mathbf{D}^* and number of non zero elements in the columns of \mathbf{A}^* , and the noise level, Gribonval et al. [2014] show that the cost function defined in Eq. (1.25) has a local minimum around $(\mathbf{D}^*, \mathbf{A}^*)$ with high probability, as long as the number of signals n is greater than a constant times mp^3 . In the noiseless case, earlier results have been also obtained Gribonval and Schnass [2010], Geng et al. [2011], and recently it has been shown that under additional assumptions, a good initializer could be found so that the previous type of local consistency results can be applied [Agarwal et al., 2013].

Finally, recent algorithms have emerged in the theoretical science community, which are not explicitly based on optimization [see, *e.g.*, Spielman et al., 2013, Recht et al., 2012, Arora et al., 2014]. These come with global convergence guarantees (with additional assumptions regarding the signals), but their empirical performance on concrete signal and image processing problems have not yet been demonstrated.

2

Discovering the Structure of Natural Images

Dictionary learning was first introduced by Olshausen and Field [1996, 1997] as an unsupervised learning technique for discovering and visualizing the underlying structure of natural image patches. The results were found impressive by the scientific community, and dictionary learning gained early success before finding numerous applications in image processing [Elad and Aharon, 2006, Mairal et al., 2008c, 2009, Yang et al., 2010a]. Without making any a priori assumption about the data except a parsimony principle, the method is able to produce dictionary elements that resemble Gabor wavelets—that is, spatially localized oriented basis functions [Gabor, 1946, Daugman, 1985], as illustrated for example in Figure 2.1.

Surprisingly, the goal of automatically learning local structures in natural images was neither originally achieved in image processing, nor in computer vision. The original motivation of Olshausen and Field [1996, 1997] was in fact to establish a relation between the statistical structure of natural images and the properties of neurons from area V1 of the mammalian visual cortex.¹ The emergence of Gabor-like patterns

¹Neuroscientists have identified several areas in the human visual cortex. Area V1 is considered to be at the earliest stage of visual processing.

from natural images was considered a significant result by neuroscientists. To better understand this fact, we will briefly say a few words about the importance of Gabor models within experimental studies of the visual cortex.

Since the pioneer work of Hubel and Wiesel [1968], it is known that some visual neurons are responding to particular image features, such as oriented edges. Specifically, displaying oriented bars at a particular location in a visual field’s subject may elicit neuronal activity in some cells. Later, Daugman [1985] demonstrated that fitting a linear model to neuronal responses given a visual stimuli may produce filters that can be well approximated by a two-dimensional Gabor function. Models based on such filters have been subsequently widely used, and are still present in state-of-the-art predictive models of the neuronal activity in V1 [Kay et al., 2008, Nishimoto et al., 2011]. By showing that Gabor-like visual patterns could be automatically obtained from the statistics of natural images, Olshausen and Field [1996, 1997] provided support for the classical Gabor model for V1 cells, even though such an intriguing phenomenon does not constitute a strong evidence that the model matches real biological processes. It is indeed commonly admitted that little is known about the early visual cortex [Olshausen and Field, 2005, Carandini et al., 2005] despite intensive studies and impressive results achieved by predictive models.

In this chapter, we show how to use dictionary learning for discovering latent structures in natural images, and we discuss other unsupervised learning techniques. We present them under the unified point of view of *matrix factorization*, which makes explicit links between different learning methods. Given a training set $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ of signals—here, natural image patches—represented by vectors \mathbf{x}_i in \mathbb{R}^m , we wish to find an approximation $\mathbf{D}\boldsymbol{\alpha}_i = \sum_{j=1}^p \alpha_i[j] \mathbf{d}_j$ for each signal \mathbf{x}_i , where $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ is a matrix whose columns are called “dictionary elements”, and the vectors $\boldsymbol{\alpha}_i$ are decomposition coefficients. In other words, we wish to find a factorization

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

where the matrix $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n]$ in $\mathbb{R}^{p \times n}$ carries the coefficients $\boldsymbol{\alpha}_i$. Many unsupervised learning techniques can be cast as ma-

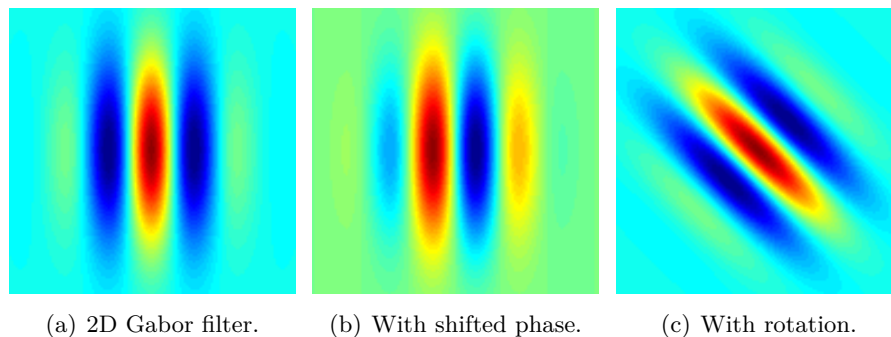


Figure 2.1: Three examples of two-dimensional Gabor filters. Red represents positive values and blue negative ones. The first two filters correspond to the function $g(x, y) = e^{-x^2/(2\sigma_x^2) - y^2/(2\sigma_y^2)} \cos(\omega x + \varphi)$, with $\varphi = 0$ for Figure (a) and $\varphi = \pi/2$ for Figure (b). Figure (c) was obtained by rotating Figure (a)—that is, replacing x and y in the function g by $x' = \cos(\theta)x + \sin(\theta)y$ and $y' = \cos(\theta)x - \sin(\theta)y$. The figure is best seen in color on a computer screen.

trix factorization problems; this includes dictionary learning, principal component analysis (PCA), clustering or vector quantization [see Nasrabadi and King, 1988, Gersho and Gray, 1992], non-negative matrix factorization (NMF) [Paatero and Tapper, 1994, Lee and Seung, 1999], archetypal analysis [Cutler and Breiman, 1994], or independent component analysis (ICA) [Hérault et al., 1985, Bell and Sejnowski, 1995, 1997, Hyvärinen et al., 2004]. These methods essentially differ in a priori assumptions that are made on \mathbf{D} and \mathbf{A} (sparse, low-rank, orthogonal, structured), and in the way the quality of the approximation $\mathbf{X} \approx \mathbf{DA}$ is measured. We will see in the rest of this chapter how these criteria influence the results obtained on natural image patches.

2.1 Pre-processing

When manipulating data, it is often important to choose an appropriate pre-processing scheme. There are several reasons for changing the way data looks like before using an unsupervised learning algorithm. For example, one may wish to reduce the amount of noise, make the data invariant to some transformation, or remove a confounding (un-

wanted) factor. Unfortunately, the literature is not clear about which pre-processing step should be applied given a specific problem at hand. In this section, we describe common practices when dealing with natural image patches, such as centering, contrast normalization, or whitening. We study the effect of these procedures on the image itself, and provide visual interpretation that should help the practitioner make his choice.

Centering. The most basic pre-processing scheme consists of removing the mean intensity from every patch. Such an approach is called “centering the data” in statistics and machine learning. Borrowing some terminology from electronics, it is also called “removing the DC component” in the signal processing literature. Formally, it means applying the following update to every signal \mathbf{x}_i :

$$\mathbf{x}_i \leftarrow \mathbf{x}_i - \left(\frac{1}{m} \sum_{j=1}^m \mathbf{x}_i[j] \right) \mathbf{1}_m,$$

where $\mathbf{1}_m$ is the vector of size m whose entries are all ones. It can also be interpreted as performing a left multiplication on the matrix \mathbf{X} :

$$\mathbf{X} \leftarrow \left(\mathbf{I} - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^\top \right) \mathbf{X}.$$

When analyzing the structure of natural patches, one is often interested in retrieving geometrical visual patterns such as edges, which is a concept related to intensity variations and thus invariant to the mean intensity. Therefore, centering can be safely used; it makes the data invariant to the mean intensity, and the learned structures are expected to have zero mean as well. The effect of such a pre-processing scheme can be visualized in Figure 2.2(b) for the image *man*. We extract all overlapping patches from the image, remove their average pixel value, and recombine the centered patches into a new image, which is displayed in Figure 2.2(b). Because of the overlap, every pixel is contained in several patches; the pixel values in the reconstructed image are thus obtained by averaging their counterparts from the centered patches. The procedure is in fact equivalent to applying a high-pass filter to the

original image. As a result, the geometrical structures are still present in the centered image, but some low frequencies have disappeared.

In practice, centering the data has been found useful in concrete applications, such as image denoising with dictionary learning [Elad and Aharon, 2006]. First removing the mean component before learning the patches structure, and then adding the mean back after processing, leads to substantially better results than working with the raw patches directly.

Variance - contrast normalization. After centering, another common practice is to normalize the signals to have unit ℓ_2 -norm, or equivalently unit variance. In computer vision, the terminology of “contrast normalization” is also often used. One motivation for such a pre-processing step is to make different regions of an image more homogeneous and to gain invariance to illumination changes. Whereas this is probably not useful for image reconstruction problems, it is a key component of many visual recognition architectures [Pinto et al., 2008, Jarrett et al., 2009]. However, contrast normalization should be applied with care: the naive update $\mathbf{x}_i \leftarrow \mathbf{x}_i / \|\mathbf{x}_i\|_2$ is likely to provide poor results for patches from uniform areas—*e.g.*, from the sky in natural scenes. After centering, such patches do not contain any information about the scene anymore but carry residual noise, which will be greatly amplified by the naive ℓ_2 -normalization. A simple way of fixing that issue consists in providing a different treatment for patches with a small norm—say, smaller than a parameter η —and use the following update instead:

$$\mathbf{x}_i \leftarrow \frac{1}{\max(\|\mathbf{x}_i\|_2, \eta)} \mathbf{x}_i.$$

We present in Figure 2.2(c) the effect of contrast normalization where η is chosen to be 0.2 times the mean value of $\|\mathbf{x}_i\|_2$ in the image, following the same patch recombination scheme as in the previous paragraph. Compared to Figure 2.2(c), areas with small intensity variations have been amplified, making the image more homogeneous. Contrast normalization can also be applied after a whitening step, which we now present.

Whitening and dimensionality reduction. Centering consists of removing from data the first-order statistics for every patch. Some early work dealing with natural images patches goes a step further by removing global second-order statistics, a procedure called *whitening* [Olshausen and Field, 1996, 1997, Bell and Sejnowski, 1997, Hyvärinen et al., 2009].

Let us consider a random variable \mathbf{x} representing centered patches uniformly sampled at random from natural images. Even though the centering step is performed at the patch level, the pixel values across centered natural images patches have zero mean $\boldsymbol{\mu} \triangleq \mathbb{E}[\mathbf{x}] = 0$. The different pixels of a patch are indeed identically distributed, and thus the vector $\boldsymbol{\mu}$ is constant; it is then necessarily equal to the constant zero vector since it has zero mean after centering. Therefore, the covariance matrix of \mathbf{x} is the quantity $\boldsymbol{\Sigma} \triangleq \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$. Whitening consists of finding a data transformation such that $\boldsymbol{\Sigma}$ becomes close to the identity matrix. As a result, the pixel values within a whitened patch are uncorrelated.

In practice, whitening is performed by computing the eigenvalue decomposition of the sample covariance matrix $(1/n) \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{U}\mathbf{S}^2\mathbf{U}^\top$, where \mathbf{U} is an orthogonal matrix in $\mathbb{R}^{m \times m}$, and \mathbf{S} in $\mathbb{R}^{m \times m}$ is diagonal with non-negative entries. Since the sample covariance matrix is positive semi-definite, the eigenvalues are non-negative and $\mathbf{S} = \text{diag}(s_1, \dots, s_m)$ carries in fact the *singular values* of the matrix $(1/\sqrt{n})\mathbf{X}$. Then, whitening consists of applying the following update to every patch \mathbf{x}_i :

$$\mathbf{x}_i \leftarrow \mathbf{U}\mathbf{S}^\dagger\mathbf{U}^\top \mathbf{x}_i,$$

where $\mathbf{S}^\dagger = \text{diag}(s_1^\dagger, \dots, s_m^\dagger)$ with $s_j^\dagger = 1/s_j$ if $|s_j| > \varepsilon$ and 0 otherwise, where ε is a small threshold to prevent arbitrarily large entries in \mathbf{S}^\dagger . It is then easy to show that the sample covariance matrix after whitening is diagonal with $k = \text{rank}(\boldsymbol{\Sigma})$ entries equal to 1. Choosing ε can be interpreted as controlling the dimensionality reduction effect of the whitening procedure. According to Hyvärinen et al. [2009], such a step is important for independent component analysis. The experiments presented in this chapter do not include dimensionality reduction for dictionary learning with grayscale image patches—that is, we choose $\varepsilon = 0$, since the resulting whitened images do not seem to suffer

from any visual artifact. The case of color patches is slightly different and is discussed in the next paragraph.

Similarly as for centering and contrast normalization, we visualize the effect of whitening on Figure 2.2(d). Removing the spatial correlation essentially results in sharper edges. As a matter of fact, whitening can be shown to amplify high frequencies and reduce low ones [Hyvärinen et al., 2009]. The reason is related to the nature of the principal components of natural images that form the columns of the matrix \mathbf{U} , which will be discussed in Section 2.2. In practice, whitening is often not used in image restoration applications based on dictionary learning since it modifies the nature of the noise.

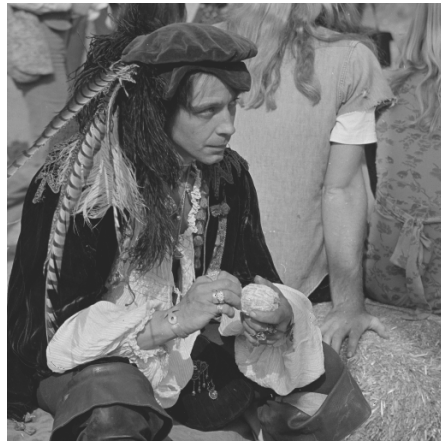
Pre-processing color patches. Finally, we conclude this section on pre-processing by discussing the treatment of color image patches. We assume that images are encoded in the RGB space, and that patches are obtained by concatenating information from the three color channels. In other words, $l \times l$ image patches are represented by vectors \mathbf{x}_i of size $m = 3 \times l \times l$.

The main question about color processing is probably whether or not the RGB color space should be used [Pratt, 1971, Faugeras, 1979, Sharma and Trussell, 1997]. The choice of the three channels R, G, and B directly originates from our first understanding of the nature of color. After discovering the existence of the color spectrum, Newton [1675] already suggests that the rays of light “impinging” on the retina will

“affect the sense with various colours, according to their bigness and mixture; the biggest with the strongest colours, reds and yellow; the least with the weakest, blue and violets; the middle with green”.

Building upon Newton’s findings, it appears that Young [1845] was the first to introduce the concept of trichromatic vision. According to Maxwell [1860]:

“Young appears to have originated the theory, that the three elements of colour are determined as much by the constitution of the sense of sight as by anything external to us. He



(a) Without pre-processing.



(b) After centering.

(c) After centering and ℓ_2 -normalization.

(d) After whitening.

Figure 2.2: The effect of various pre-processing procedures on the image man. For the figures (b), (c) and (d), the value 0 is represented by a gray pixel. Since they contain negative values, pre-processed images are shifted and rescaled for visualization purposes.

conceives that three different sensations may be excited by light. . . He conjectures that these primary sensations correspond to red, green, and violet.”

Finally, von Helmholtz [1852], Grassmann [1854], and Maxwell [1860] proposed rigorous rules of color compositions, paving the way to the modern treatment of color in vision processing [see Forsyth and Ponce, 2012, Chapter 3]. More than a century after these early discoveries, the existence in the eye of three types of biological photoreceptors—respectively corresponding to red, green, and blue wavelengths—is established [Nathans et al., 1986], and the RGB color space is widely used in electronic devices.

The question of changing the color space can be rephrased as follows: should we apply a linear or non-linear transformation f to each RGB pixel value: $(u, v, w) = f(r, g, b)$ and work in the resulting space? One of the early motivation for studying the representation of color was to reduce the bandwidth required by RGB over communication channels. RGB components are indeed highly correlated; thus, removing the redundant information between the different channels allows more efficient coding for transmitting information [Pratt, 1971]. Another motivation is that the Euclidean distance on RGB is known to be a poor estimation of the perceptual distance by humans. Therefore, a large variety of color spaces have been designed, such as CIE Lab, CIE XYZ, YIQ, or YCrBr [see Sharma and Trussell, 1997], which are partially improving upon RGB regarding the decorrelation of the transformed channels and human perception. However, it remains unclear whether these spaces should be used or not in practice, and the optimal choice depends on the task and algorithm. For instance, for color image denoising, Takeda et al. [2007] and Dabov et al. [2007b] choose the YCrBr space, whereas Buades et al. [2005], Mairal et al. [2008c] and Chatterjee and Milanfar [2012] successfully use RGB. Changing the color space modifies indeed the noise structure as a side effect, which may be problematic for some methods.

The other pre-processing steps that we have presented for grayscale images can also be applied to color image patches. The first difference with the monochromatic setting lies in the centering scheme when work-

ing with RGB. One motivation for centering grayscale image patches is to learn geometric structures that are invariant to the mean intensity. For color images, the quantity $\frac{1}{m} \sum_{j=1}^m \mathbf{x}_i[j]$ unfortunately lacks of interpretation, and removing it from a patch does not seem to achieve any interesting property. Instead, we may be looking for geometrical structures in images that are invariant to color patterns, and thus, one may remove the mean color from the patch. In other words, one should *center every R, G, B channel independently*. As a result, the centering scheme becomes also invariant to any linear transformation of the color space, which can be a desirable property. The effect of such a preprocessing can be visualized in Figure 2.3(b) for the image kodim07 from the Kodak PhotoCD dataset.² Note that centered images have negative values and thus images are shifted and rescaled for visualization purposes. Regions with relatively uniform colors appear grayish, and thus mostly contain geometrical content; some other areas contain transitions between a color—say, represented by RGB values (r, g, b) —and its opponent one—represented by $(-r, -g, -b)$. The colored halos around the pink flowers and the green leaves are thus due to this color/opponent color transitions.

In the experiments of this paper about dictionary learning, whitening comprises a light dimensionality reduction step, where we threshold to zero the smallest singular values of \mathbf{X} (the entries on the diagonal of \mathbf{S}). More precisely, we keep the k largest singular values s_1, \dots, s_k such that $\sum_{i=1}^k s_i^2 \geq 0.995 \sum_{i=1}^m s_i^2$; in other words, we keep 99.5% of the explained variance of the data. Contrast normalization can then be applied without particular attention to the nature of the patch. The effects of whitening and contrast normalization are displayed on Figures 2.3(c) and 2.3(d), respectively. Similar conclusions can be drawn for color image patches as in the monochromatic case.

2.2 Principal component analysis

Principal component analysis (PCA), also known as the Karhunen-Loève or Hotelling transform [Hotelling, 1933], is probably the most

²available here: <http://r0k.us/graphics/kodak/>.

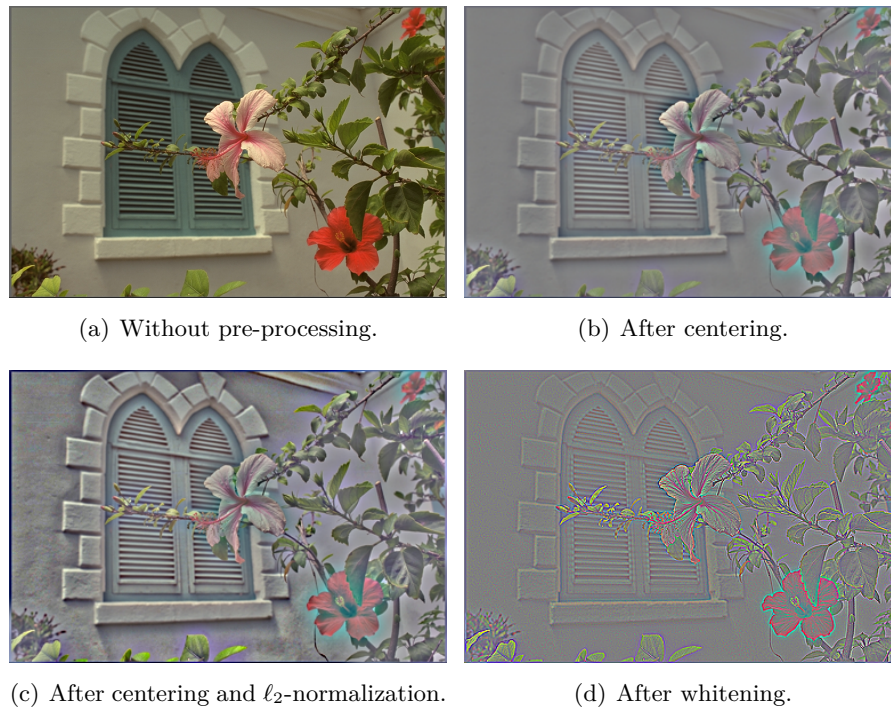


Figure 2.3: The effect of various pre-processing procedures on the image kodim07. Pre-processed images are shifted and rescaled for visualization purposes, since they contain negative values. The figure is best seen in color on a computer screen.

widely used unsupervised data analysis technique. Even though it is often presented as an iterative process finding orthogonal directions maximizing variance in the data, it can be cast as a low-rank matrix factorization problem:

$$\min_{\mathbf{U} \in \mathbb{R}^{m \times k}, \mathbf{V} \in \mathbb{R}^{n \times k}} \left\| \mathbf{X} - \mathbf{U}\mathbf{V}^\top \right\|_F^2 \quad \text{s.t.} \quad \mathbf{U}^\top \mathbf{U} = \mathbf{I}_k,$$

where k is the number of principal components we wish to obtain, and \mathbf{I}_k is the identity matrix in $\mathbb{R}^{k \times k}$. We also assume that the rows of the matrix \mathbf{X} have zero mean. As a consequence of the theorem of Eckart and Young [1936], the matrix \mathbf{U} contains the principal components of \mathbf{X} corresponding to the k largest singular values. In Figure 2.4(b), we visualize the principal components of $n = 400\,000$ natural image patches of size 16×16 pixels, ordered by largest to smallest singular value (from left to right, then top to bottom). The components resemble Fourier basis—that is, product of sinusoids with different frequencies and phases, similarly to the discrete cosine transform (DCT) dictionary [Ahmed et al., 1974] presented in Figure 2.4(a).

However, there is a good reason for obtaining sinusoids that is simply related to a property of translation invariance, *meaning that the patterns observed in Figure 2.4(b) are unrelated to the underlying structure of natural images*. This fact is well known and has been pointed earlier by others in the literature [see, e.g., Bossomaier and Snyder, 1986, Field, 1987, Simoncelli and Olshausen, 2001, Hyvärinen et al., 2009]. Second-order statistics of natural images patches are commonly assumed to be invariant by translation, such that the correlation of pixel values at locations $z_1 = (k_1, l_1)$ and $z_2 = (k_2, l_2)$ only depends on the displacement $z_1 - z_2$ [Simoncelli and Olshausen, 2001]. This is particularly true for patches that are extracted from larger images at any arbitrary position, and whose distribution is exactly translation invariant. Such signals are called “stationary” and their principal components (equivalently the eigenvectors of the covariance matrix) are often considered to be well approximated by the Discrete Fourier Transform (DFT), as noted by Pearl [1973], leading to sinusoidal principal components.

Nevertheless, the relation between principal components and si-

nusoids only holds rigorously in an asymptotic regime. For instance, consider the case of an infinite one-dimensional signal with covariance $\Sigma[k, l] = \sigma(k - l)$ for positions k and l , where σ is an even function. Then, for all frequency ω and phase φ ,

$$\sum_l \Sigma(k, l) e^{i(\omega l + \varphi)} = \sum_l \sigma(l - k) e^{i(\omega l + \varphi)} = \left(\sum_{l'} \sigma(l') e^{i\omega l'} \right) e^{i(\omega k + \varphi)},$$

where i denotes the imaginary unit, and the sums are over all integers. Since the function σ is even, the infinite sum $(\sum_{l'} \sigma(l') e^{i\omega l'})$ is real, and the signals $[\sin(\omega k + \varphi)]_{k \in \mathbb{Z}}$ are all eigenvectors of the covariance operator Σ . Equivalently, they are principal components of the input data. The case of two-dimensional signals can be treated similarly but it involves heavier notation. Drawing conclusions about the eigenvectors of a finite covariance matrix computed on natural image patches is nevertheless subject to discussion, and understanding the quality of the approximation of the principal components by the DFT in the finite regime is non-trivial [Pearl, 1973].

From an experimental point of view, the fact that the structure of natural images has nothing to do with the patterns displayed in Figure 2.4(b) is easy to confirm. In Figure 2.5, we present principal components computed on all overlapping patches from the image tiger. Even though this image is sketched by hand—and thus, is not natural—we recover sine waves as in Figure 2.4(b).

2.3 Clustering or vector quantization

Clustering techniques have been used for a long time on natural image patches for compression and communication purposes under the name of “vector quantization” [Nasrabadi and King, 1988, Gersho and Gray, 1992]. The goal is to find p clusters in the data, by minimizing the following objective:

$$\min_{\substack{\mathbf{D} \in \mathbb{R}^{m \times p} \\ \forall i, l_i \in \{1, \dots, p\}}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{d}_{l_i}\|_2^2, \quad (2.1)$$

where the columns of $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ are called “centroids” and l_i is the index of the cluster associated to the data point \mathbf{x}_i . The al-

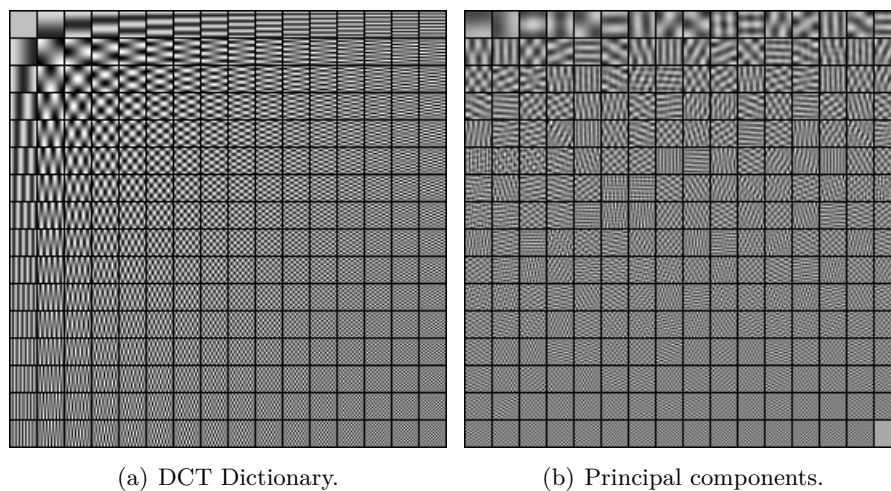


Figure 2.4: On the right, we visualize the principal components of 400 000 randomly sampled natural image patches of size 16×16 , ordered by decreasing variance, from top to bottom and left to right. On the left, we display a discrete cosine transform (DCT) dictionary. Principal components resemble DCT dictionary elements.

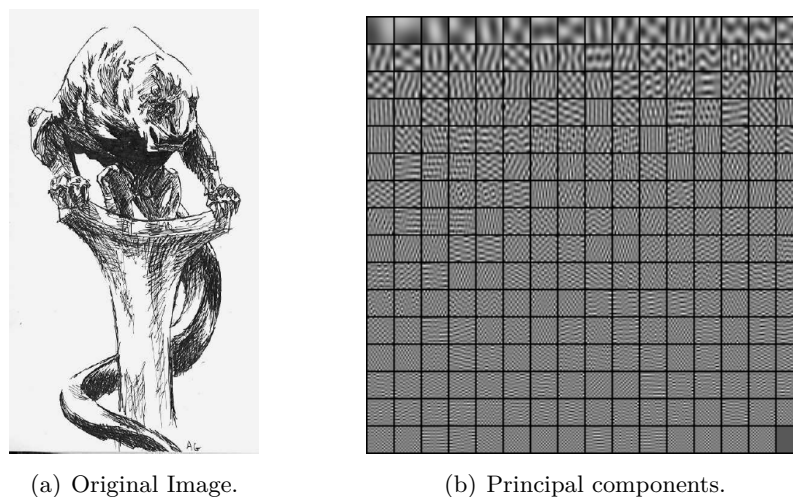


Figure 2.5: Visualization of the principal components of all overlapping patches from the image tiger. Even though the image is not natural, its principal components are similar to the ones of Figure 2.4(b).

gorithm K-means [see Hastie et al., 2009] approximately optimizes the non-convex objective (2.1) by alternatively performing exact minimization with respect to the labels l_i with \mathbf{D} fixed, and with respect to \mathbf{D} with the labels fixed.

To make the link between clustering and matrix factorization, it is also possible to reformulate (2.1) as follows

$$\min_{\substack{\mathbf{D} \in \mathbb{R}^{m \times p} \\ \mathbf{A} \in \{0,1\}^{p \times n}}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \quad \text{s.t.} \quad \forall i, \sum_{j=1}^p \boldsymbol{\alpha}_i[j] = 1,$$

where $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n]$ carries binary vectors that sum to one, and clustering can be subsequently seen as a matrix factorization problem:

$$\min_{\substack{\mathbf{D} \in \mathbb{R}^{m \times p} \\ \mathbf{A} \in \{0,1\}^{p \times n}}} \frac{1}{2n} \|\mathbf{X} - \mathbf{D}\mathbf{A}\|_F^2 \quad \text{s.t.} \quad \forall i, \sum_{j=1}^p \boldsymbol{\alpha}_i[j] = 1.$$

Then, the algorithm K-means is performing alternate minimization between \mathbf{A} and \mathbf{D} , each step decreasing the value of the objective.

We visualize clustering results in Figure 2.6 after applying 250 iterations of the K-means algorithm on $n = 400\,000$ natural image patches and choosing $p = 256$ centroids. Without whitening, K-means produces mostly low-frequency patterns, which is not surprising since most of the energy (ℓ_2 -norm) of images is concentrated in low frequencies. It is indeed known that the spectral power of natural images obtained in the Fourier domain typically decreases according to the power law $1/f^2$ for a frequency f [see Simoncelli and Olshausen, 2001]. After whitening, high-frequency “Gabor-like” patterns and checkerboard patterns emerge from data. It is thus interesting to see that with an appropriate pre-processing procedure, simple unsupervised learning techniques such as K-means are able to discover Gabor features with different orientations, frequencies, and positions within patches. We note that the localized checkerboard patterns might be an artifact of the whitening procedure without dimensionality reduction, which amplifies such patterns appearing among the last principal components presented in Figure 2.4(b).

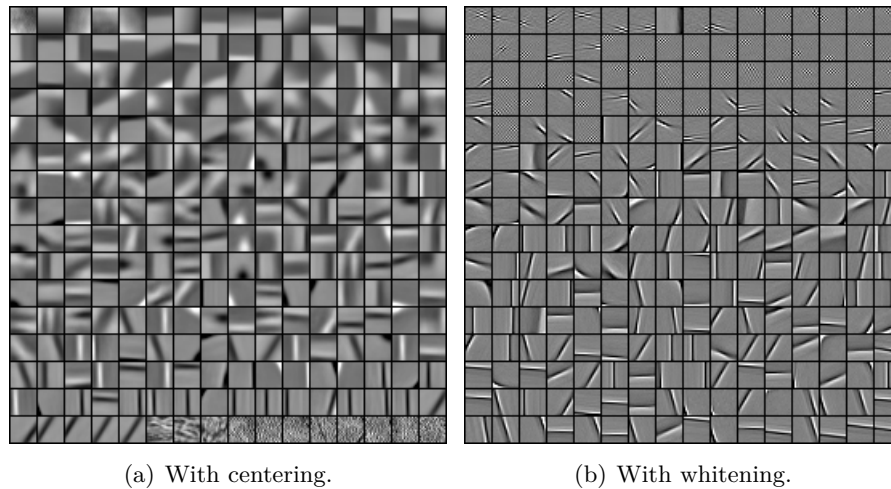


Figure 2.6: Visualization of $p = 256$ centroids computed with the algorithm K-means on $n = 400\,000$ image patches of size $m = 16 \times 16$ pixels. We compare the results with centering (left), and whitening (right). In both cases, we also apply a contrast normalization step. Centroids are ordered from most to least used (from left to right, then top to bottom).

2.4 Dictionary learning

We are now interested in the dictionary learning formulation originally introduced by Olshausen and Field [1996, 1997] and its application to natural image patches. We consider the corresponding matrix factorization formulation (1.22), which we recall here

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{D}\mathbf{A}\|_{\text{F}}^2 + \lambda \Psi(\mathbf{A}), \quad (2.2)$$

where $\Psi(\mathbf{A}) = \sum_{i=1}^n \psi(\alpha_i)$ and \mathcal{C} is the set of matrix whose columns have a Euclidean norm smaller than one.

In Figure 2.7, we present visual results obtained with such a formulation where ψ is the ℓ_1 -norm. As in the previous section about clustering, we use $n = 400\,000$ natural image patches, either gray, or extracted from color images in RGB. We center the patches, according to the procedure presented in Section 2.1, and rescale the matrix \mathbf{X} such that its columns have unit norm on average. Then, we learn a dictionary \mathbf{D} with $p = 256$ dictionary elements and set the regularization parameter λ to 0.1. We use the online dictionary learning algorithm of Mairal et al. [2010a] available in the software SPAMS, making 10 passes over the data. The dictionaries obtained for gray and RGB patches and with two different preprocessing steps are displayed in Figure 2.7. Note that other dictionary learning approaches can be used for this task and providing similar results. This is for instance the case of the K-SVD algorithm of Aharon et al. [2006], which we present in Section 5.5.

When processing grayscale image patches, a simple centering step produces dictionary elements with both low-frequency elements, and high-frequency, localized, Gabor-like patterns, as shown in Figure 2.7(a). The results after a whitening step are displayed in Figure 2.7(c). Whitening corresponds to applying a high-pass filter to the original images, and we obtain exclusively Gabor-like features as in the early work of Olshausen and Field [1996, 1997]. The case of color image patches presented in Figures 2.7(b) and 2.7(d) is also interesting. After centering, about 230 dictionary elements out of 256 are grayscale, or are almost grayscale, which is consistent with the greyish appearance

of images after centering (see Section 2.1). An intriguing phenomenon concerns the remaining colored dictionary elements. Those exhibit low-frequency patterns with two opponent colors, blue and yellow for 10 of them, green and red/purple for 16 of them. Interestingly, such opponent color patterns are typical characteristics of neurons' receptive fields described in the neuroscience literature [Livingstone and Hubel, 1984, Ts'o and Gilbert, 1988].

2.5 Structured dictionary learning

Because the ℓ_1 -norm cannot model interactions between dictionary elements, it is natural to consider extensions of dictionary learning where a particular structure is taken into account. The concept of *structured sparsity* presented in Section 1.3 is a natural tool for this; it consists of replacing the ℓ_1 -norm by a more complex sparsity-inducing penalty.

Hierarchical dictionary learning. A first extension has been investigated by Jenatton et al. [2010a, 2011b], when a pre-defined hierarchical structure is assumed to exist among the p dictionary elements. The formulation (2.2) is considered where the function Ψ is the hierarchical Group-Lasso penalty of Zhao et al. [2009]. A tree is given, *e.g.*, is defined by the user, and one dictionary element is associated to every node of the tree. The penalty constructs a group structure \mathcal{G} of subsets of $\{1, \dots, p\}$, each group containing one node and all its descendants in the tree. An example is illustrated in Figure 1.12. The resulting penalty Ψ is then defined as

$$\Psi(\mathbf{A}) = \sum_{i=1}^n \sum_{g \in \mathcal{G}} \|\alpha_i[g]\|_q, \quad (2.3)$$

where $\|\cdot\|_q$ can either represent the ℓ_∞ - or the ℓ_2 -norm. As explained in Section 1.3, an effect of the penalty is that a dictionary element can be used in the decomposition of a patch only if its parent in the tree is also used. We present some visual results for $q = \infty$ in Figure 2.8, after centering the natural image patches, for some manually tuned regularization parameter λ , and two different tree structures. Dictionary

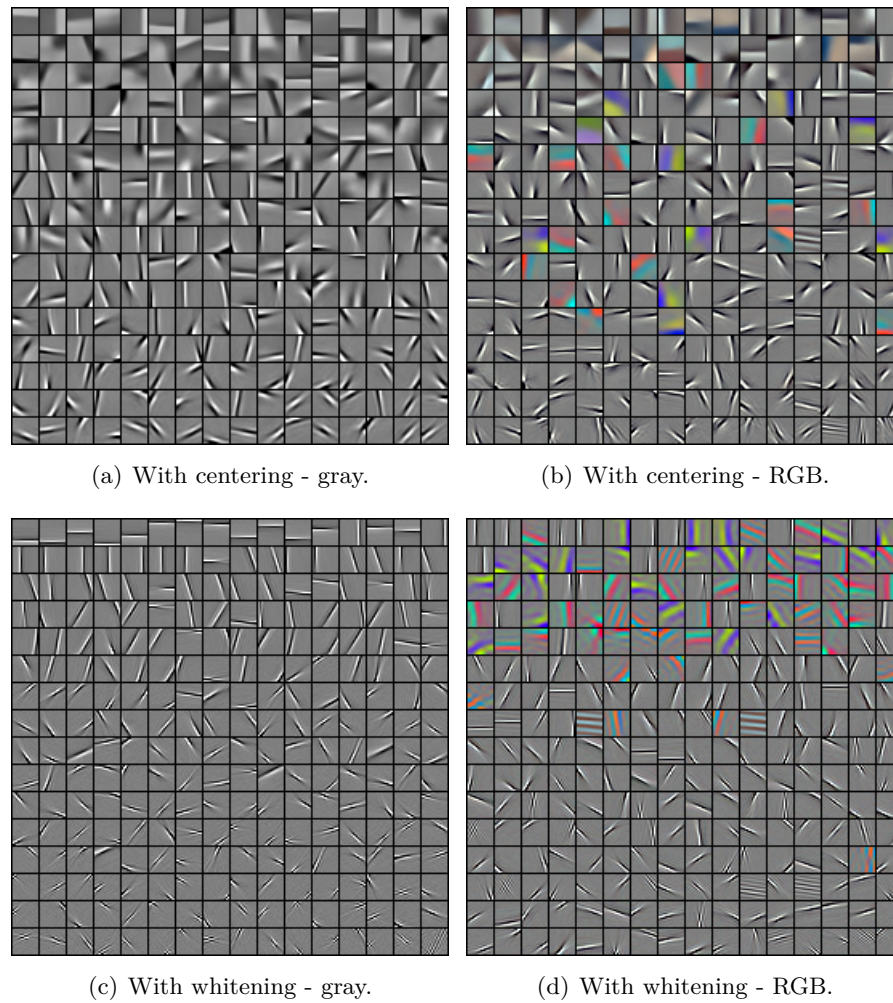


Figure 2.7: Dictionaries obtained with the formulation (2.2), with two different pre-processing procedures, and with gray or RGB patches. The dictionary elements are ordered from most to least used (from left to right, then top to bottom). Best seen in color on a computer screen.

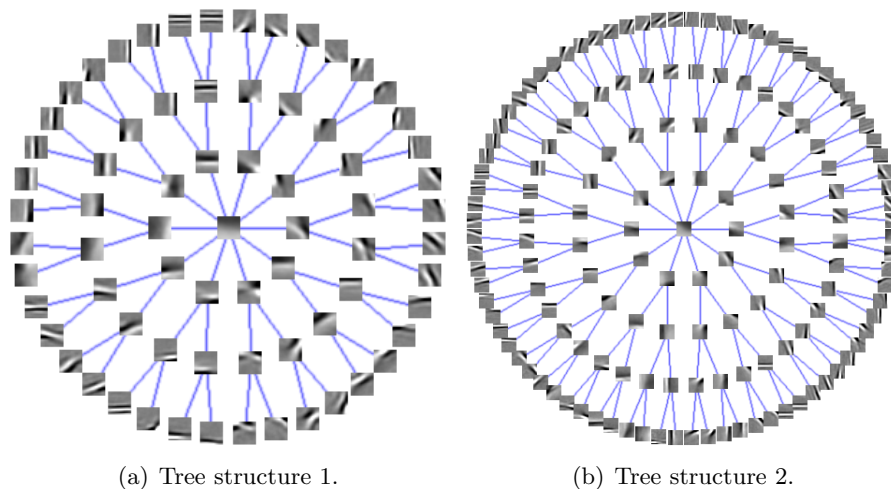


Figure 2.8: Dictionaries obtained with the formulation (2.2) and a hierarchical structured-sparsity penalty Ψ for natural image patches of size 16×16 pixels. The tree structures are pre-defined, and the dictionary elements naturally organize themselves in the tree during learning. For each tree, the root is represented in the middle of the figure. (a): the tree is of depth 4 and the branching factors at depths 1, 2, 3 are respectively 10, 2, 2; (b): the tree is slightly more complex; its depth is 5 and the branching factors at depths 1, 2, 3, 4 are respectively 10, 2, 2, 2. Figures borrowed from Jenatton et al. [2010a, 2011b]. Best seen by zooming on a computer screen.

elements naturally organize themselves in the tree, often with low frequencies near the root of the tree, and high frequencies near the leaves. We also observe strong correlations between each parent node and their children in the tree, where children often look like their parent, but with higher frequencies and with minor variations.

Topographic dictionary learning. A two-dimensional grid structure has also been used for learning dictionaries [Kavukcuoglu et al., 2009, Mairal et al., 2011], which was directly inspired from the topographic independent component analysis formulation of Hyvärinen et al. [2001]. The principle is similar to the hierarchical case—that is, we can use a penalty Ψ as in (2.3), but the set of groups \mathcal{G} takes into account a different structure. We assume here that the dictionary elements are organized on a grid, such that we can define neighborhood relations

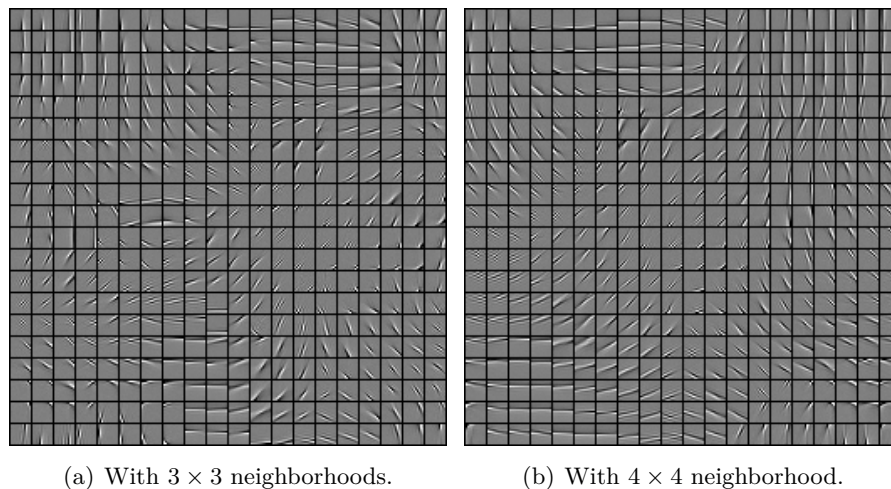


Figure 2.9: Dictionaries obtained with the formulation (2.2) and a structured sparsity penalty Ψ inducing a grid structure in the dictionary. The dictionaries are computed on whitened natural image patches of size 12×12 pixels for two different group structures. Figure borrowed from Mairal et al. [2011]. Best seen on a computer screen.

between them.

For instance, we can organize the p dictionary elements on a $\sqrt{p} \times \sqrt{p}$ grid, and consider p overlapping groups that are 3×3 or 4×4 spatial neighborhoods on the grid (to avoid boundary effects, we assume the grid to be cyclic). We represent some results obtained with such a formulation in Figure 2.9, where a function Ψ is defined in (2.3) with $q = 2$. The regularization parameter λ is also manually tuned. The dictionary elements automatically organize themselves on the grid, and exhibit some intriguing spatial smoothness. Another formulation achieving a similar effect was also proposed by Garrigues and Olshausen [2010], by mixing sparse coding with a probabilistic model involving latent variables for groups of neighbor variables. Another approach was also proposed by Gregor et al. [2011], which models inhibition effects between dictionary elements.

2.6 Other matrix factorization methods

We have focused our study of learning methods for natural image patches on sparse coding and clustering techniques so far. Other matrix factorization formulations have shown to be also effective for this task. In particular, independent component analysis (ICA) was applied successfully to image patches around the same time as dictionary learning by Bell and Sejnowski [1997], who also obtained Gabor-like patterns from whitened data.

Independent component analysis. A significantly different point of view than dictionary learning consists of factorizing *whitened* data \mathbf{X} as a product $\mathbf{D}\mathbf{A}$, where the sparsity principle is replaced by an assumption of statistical independence. Therefore, the ICA approach requires by nature a probabilistic interpretation of the data. We assume that \mathbf{x} in \mathbb{R}^m is a random variable representing a signal—here, a natural image patch—and the columns of $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ are random realizations of the variable \mathbf{x} —that is, a set of natural image patches selected uniformly at random. The principle of ICA assumes that there exists a factorization $\mathbf{x} = \mathbf{D}\boldsymbol{\alpha}$, where \mathbf{D} is an orthogonal matrix, and $\boldsymbol{\alpha} = \mathbf{D}^\top \mathbf{x}$ is a random vector whose entries are statistically independent [see Bell and Sejnowski, 1997, Hyvärinen et al., 2004, 2009]. Then, the columns of \mathbf{D} are called independent components.

Given data, the concept of “statistical independence” alone does not immediately yield a precise cost function to optimize, and thus numerous variants of ICA exist in the literature. Ultimately, most variants of ICA try to approximate the statistical independence assumption, by finding an orthogonal transformation \mathbf{D}^\top such that the probability density $p(\boldsymbol{\alpha})$, assuming it exists, factorizes into the product of its marginals $\prod_{j=1}^p p(\boldsymbol{\alpha}[j])$. A popular cost function to compare probability distributions is the Kullback-Leibler distance [see Cover and Thomas, 2006], denoted by KL , and a natural goal of ICA is to minimize over \mathbf{D}

$$KL \left(p(\boldsymbol{\alpha}), \prod_{j=1}^p p(\boldsymbol{\alpha}[j]) \right) \triangleq \int_{\mathbb{R}^p} p(\boldsymbol{\alpha}) \log \left(\frac{p(\boldsymbol{\alpha})}{\prod_{j=1}^p p(\boldsymbol{\alpha}[j])} \right) d\boldsymbol{\alpha},$$

which is equal to zero if and only if the $\boldsymbol{\alpha}[j]$ ’s are independent. The

quantity KL above is also called the mutual information between the variables $\alpha[j]$, and can also be simply written in terms of the entropy function H :³

$$KL \left(p(\boldsymbol{\alpha}), \prod_{j=1}^p p(\alpha[j]) \right) = \sum_{j=1}^p H(\alpha[j]) - H(\boldsymbol{\alpha}).$$

Interestingly, the entropy $H(\boldsymbol{\alpha})$ can be shown to be independent of \mathbf{D} when \mathbf{D} is orthogonal and \mathbf{x} is whitened with identity covariance, such that the goal of ICA can be expressed as the minimization of the sum of entropies

$$\sum_{j=1}^p H(\alpha[j]) = \sum_{j=1}^p H(\mathbf{d}_j^\top \mathbf{x}). \quad (2.4)$$

Unfortunately, entropy remains an abstract quantity that is not computable, and thus (2.4) is not yet an objective function that is easy to minimize. Turning (2.4) into a concrete formulation and algorithm can be achieved by following different strategies. One of them consists of parameterizing the densities $p(\mathbf{d}_j^\top \mathbf{x})$, assuming they admit a known parametric form, and minimizing (2.4) becomes equivalent to performing maximum likelihood estimation [see Hyvärinen et al., 2004]. Another approach consists of using non-parametric estimators of the entropy [Pham, 2004]. Finally, a large class of methods consists of using approximations of the entropy cost function (2.4) or encourage the distributions of the $\alpha[j]$'s to be “non-Gaussian” [Cardoso, 2003]. The non-gaussianity principle might not seem intuitive at first sight, but it is in fact rather natural when considering a minimization entropy problem. Among all probability distributions with same variance, the Gaussian ones are known to maximize entropy [Cover and Thomas, 2006]. As such, minimizing (2.4) implies that we are looking for some non-Gaussian random variables $\alpha[j]$.

Extensions where \mathbf{D} is not orthogonal and when the data \mathbf{X} is not whitened exist, but we have omitted them in this brief presentation for simplicity [see Hyvärinen et al., 2004, for a better overview].

³The entropy is defined as $H(x) \triangleq \int p(x) \log p(x) dx$ [see Cover and Thomas, 2006]. In general, $H(x)$ is not computable exactly since the probability density p is unknown.

A remarkable fact is that ICA applied to natural image patches yields very similar results to the sparse coding model of Olshausen and Field [1996, 1997], as shown by Bell and Sejnowski [1997]. In fact, an equivalence exists between the two formulations under particular asymptotic conditions [see the appendix of Olshausen and Field, 1997].

Non-negative matrix factorization. Another popular unsupervised learning technique is the non-negative matrix factorization (NMF) [Paatero and Tapper, 1994], which was shown to be able to automatically discover some interpretable features on datasets of human faces [Lee and Seung, 1999]. When the data matrix \mathbf{X} is non-negative, the method consists of finding a product \mathbf{DA} that approximates \mathbf{X} , and where each factor is also non-negative. The corresponding formulation can be written as

$$\min_{\mathbf{D} \in \mathbb{R}^{m \times p}, \mathbf{A} \in \mathbb{R}^{p \times n}} \|\mathbf{X} - \mathbf{DA}\|_{\mathbb{F}}^2 \quad \text{s.t.} \quad \mathbf{D} \geq 0 \quad \text{and} \quad \mathbf{A} \geq 0.$$

Even though this technique could be applied to natural image patches without pre-processing, it has been empirically observed that NMF does not yield any interpretable feature when applied to such data [Mairal et al., 2010a]. Note that NMF is often used with a different loss function than the square loss, in particular for audio applications [Févotte et al., 2009].

Archetypal analysis. Around the same time as dictionary learning and non-negative matrix factorization, archetypal analysis was introduced by Cutler and Breiman [1994] for discovering latent factors from high-dimensional data. The main motivation was to propose an unsupervised learning technique that is more interpretable than principal component analysis. The method seeks a factorization $\mathbf{X} = \mathbf{DA}$, with two symmetrical geometrical constraints. First, each column \mathbf{d}_j of \mathbf{D} , called “archetype” is forced to be a convex combination of a few data points. In other words, for all $j \in \{1, \dots, p\}$, there exists a vector β_j such that $\mathbf{d}_j = \mathbf{X}\beta_j$, and β_j is in the simplex Δ_n defined as

$$\Delta_n \triangleq \left\{ \beta \in \mathbb{R}^n \quad \text{s.t.} \quad \beta \geq 0 \quad \text{and} \quad \sum_{i=1}^n \beta[i] = 1 \right\}. \quad (2.5)$$

Second, each data point \mathbf{x}_i is encouraged to be close to the convex hull of the archetypes, meaning that \mathbf{x}_i should be close to a product $\mathbf{D}\boldsymbol{\alpha}_i$, where $\boldsymbol{\alpha}_i$ is in the simplex Δ_p , defined as in (2.5) when replacing n by p . The resulting formulation is the following

$$\min_{\substack{\boldsymbol{\alpha}_i \in \Delta_p \text{ for } 1 \leq i \leq n \\ \boldsymbol{\beta}_j \in \Delta_n \text{ for } 1 \leq j \leq p}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \quad \text{s.t.} \quad \forall j, \mathbf{d}_j = \mathbf{X}\boldsymbol{\beta}_j,$$

which, in the matrix factorization form, is equivalent to

$$\min_{\substack{\boldsymbol{\alpha}_i \in \Delta_p \text{ for } 1 \leq i \leq n \\ \boldsymbol{\beta}_j \in \Delta_n \text{ for } 1 \leq j \leq p}} \|\mathbf{X} - \mathbf{X}\mathbf{B}\mathbf{A}\|_F^2,$$

where $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n]$, $\mathbf{B} = [\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_p]$ and the matrix of archetypes \mathbf{D} is equal to the product $\mathbf{X}\mathbf{B}$.

Archetypal analysis is closely related to the dictionary learning formulation of Olshausen and Field [1996, 1997] since the simplicial constraint $\boldsymbol{\alpha}_i \in \Delta_p$ can be equivalently rewritten as the ℓ_1 -norm constraint $\|\boldsymbol{\alpha}_i\|_1 = 1$ associated to a non-negativity one $\boldsymbol{\alpha}_i \geq 0$. As a result, the coefficients $\boldsymbol{\alpha}_i$ are sparse in practice and archetypal analysis provides a sparse decomposition of the data points \mathbf{x}_i . The main difference with dictionary learning is the set of constraints $\mathbf{d}_j = \mathbf{X}\boldsymbol{\beta}_j$. Because the vectors $\boldsymbol{\beta}_j$ are constrained to be in the simplex Δ_n , they are encouraged to be sparse and each archetype is a convex combination of a few data points only. Such a relation between latent factors \mathbf{d}_j and the data \mathbf{X} is useful whenever interpreting \mathbf{D} is important, *e.g.*, in experimental sciences. For example, clustering techniques provide such associations between data and centroids. It is indeed common in genomics to cluster gene expression data from several individuals, and to interpret each centroid by looking for some common physiological traits among individuals of the same cluster [Eisen et al., 1998].

Archetypal analysis is also related to non-negative matrix factorization [Paatero and Tapper, 1994]. When the data \mathbf{X} is non-negative, it is easy to see that both the matrix \mathbf{A} and the matrix \mathbf{D} of archetypes are also non-negative. Unfortunately, archetypal analysis did not encounter as much success as dictionary learning or NMF, despite the fact that it provides an elegant methodology for interpreting its output.

One of the reason for this lack of popularity may be that no efficient software has been available for a long time, which may have limited its application to important scientific problems. Based upon such observations, Chen et al. [2014] have recently revisited archetypal analysis for computer vision, with the goal of bringing back this powerful unsupervised learning technique into favor. They made publicly available an efficient implementation of archetypal analysis in the SPAMS software, and demonstrated that it could perform as well as dictionary learning for some classification tasks, while offering natural mechanisms for visualizing large databases of images (see Section 4.6).

Bayesian models for dictionary learning. Early models of dictionary learning were probabilistic [see, for instance, Lewicki and Olshausen, 1999, Lewicki and Sejnowski, 2000]. Each image patch \mathbf{x} is considered to be a random variable with a normal distribution given a dictionary \mathbf{D} and a latent variable $\boldsymbol{\alpha}$. More precisely, we have the conditional law

$$p(\mathbf{x}|\mathbf{D}, \boldsymbol{\alpha}) \propto e^{-\frac{1}{2\sigma^2}\|\mathbf{x}-\mathbf{D}\boldsymbol{\alpha}\|_2^2}, \quad (2.6)$$

and the prior distribution on the coefficients $\boldsymbol{\alpha}$ is Laplacian:⁴

$$p(\boldsymbol{\alpha}) \propto e^{-\lambda\|\boldsymbol{\alpha}\|_1}. \quad (2.7)$$

With this probabilistic model, maximizing the posterior $p(\boldsymbol{\alpha}|\mathbf{x}, \mathbf{D})$ given some observation \mathbf{x} and a dictionary \mathbf{D} , or equivalently minimizing the negative log posterior, yields the Lasso formulation. Indeed, by using Bayes' rule, we have

$$\begin{aligned} -\log p(\boldsymbol{\alpha}|\mathbf{x}, \mathbf{D}) &= -\log p(\mathbf{x}|\mathbf{D}, \boldsymbol{\alpha}) - \log p(\boldsymbol{\alpha}) \\ &= \frac{1}{2\sigma^2}\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda\|\boldsymbol{\alpha}\|_1. \end{aligned} \quad (2.8)$$

When n natural image patches $\mathbf{x}_1, \dots, \mathbf{x}_n$ are observed, learning a dictionary \mathbf{D} can be achieved by (i) modeling each \mathbf{x}_i according to (2.6) with a latent variable $\boldsymbol{\alpha}_i$ associated to \mathbf{x}_i , (ii) assuming the $\boldsymbol{\alpha}_i$'s to be

⁴Note that choosing the Laplace distribution for the prior does not mean implicitly assuming that the true distribution of the latent variable $\boldsymbol{\alpha}$ is Laplace. In fact, Gribonval et al. [2012] have shown that if this was the case, the choice (2.7) would be inappropriate in the context of maximum a posteriori estimation.

statistically independent one from each other, (iii) choosing a prior distribution $p(\mathbf{D})$, *e.g.*, uniform on the set of matrices \mathcal{C} defined in (2.2), and (iv) computing the point estimate

$$\min_{\mathbf{D}, \alpha_1, \dots, \alpha_n} -\log p(\mathbf{D}, \alpha_1, \dots, \alpha_n | \mathbf{x}_1, \dots, \mathbf{x}_n).$$

By using again Bayes' rule and by using the statistical independence assumption, we recover the classical “matrix factorization” formulation presented earlier:

$$\begin{aligned} -\log p(\mathbf{D}, \mathbf{A} | \mathbf{X}) &= -\log p(\alpha_1, \dots, \alpha_n | \mathbf{X}, \mathbf{D}) - \log p(\mathbf{D}) \\ &= \sum_{i=1}^n -\log p(\alpha_i | \mathbf{x}_i, \mathbf{D}) - \log p(\mathbf{D}) \\ &= \sum_{i=1}^n \frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1 - \log p(\mathbf{D}), \end{aligned} \tag{2.9}$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and $\mathbf{A} = [\alpha_1, \dots, \alpha_n]$. Simply maximizing the posterior distribution is however not satisfactory from the point of view of a Bayesian statistician, who is usually not interested in point estimates, but in the full posterior distribution [see Bayarri and Berger, 2004].⁵ In fact, the formulation (2.9) discards any uncertainty regarding the model parameters \mathbf{D} and \mathbf{A} and would be called a “frequentist” approach in statistics: assuming that one can repeatedly draw natural image patches at random, one wishes to find a dictionary \mathbf{D} that is good on average, which is nothing else than the (non-probabilistic) empirical risk minimization point of view presented in Section 1.4.

A first step towards a Bayesian dictionary learning formulation consists of integrating instead of maximizing with respect to the latent variables α_i , as typically done in Bayesian sparse linear mod-

⁵It seems that a common misconception in the computer vision and image processing literature is to call “Bayesian” any probabilistic model that uses Bayes' rule such as (2.8). Such a terminology for maximum a posteriori (MAP) estimation is slightly misleading regarding the hundred-year-old debate among frequentists and Bayesian statisticians. In general, the latter do not assume that there exists a “true” parameter that can be obtained with MAP estimation; instead, they treat all parameters as random variables and model their uncertainty [see Bayarri and Berger, 2004].

els [Park and Casella, 2008, Seeger, 2008], and maximizing the posterior $p(\mathbf{D}|\mathbf{X})$. As a result, we need to minimize

$$\begin{aligned}
 -\log p(\mathbf{D}|\mathbf{X}) &= -\log p(\mathbf{X}|\mathbf{D}) - \log p(\mathbf{D}) \\
 &= \sum_{i=1}^n -\log p(\mathbf{x}_i|\mathbf{D}) - \log p(\mathbf{D}) \\
 &= \sum_{i=1}^n -\log \left(\int_{\boldsymbol{\alpha}_i \in \mathbb{R}^p} p(\mathbf{x}_i, \boldsymbol{\alpha}_i|\mathbf{D}) d\boldsymbol{\alpha}_i \right) - \log p(\mathbf{D}) \\
 &= \sum_{i=1}^n -\log \left(\int_{\boldsymbol{\alpha}_i \in \mathbb{R}^p} p(\mathbf{x}_i|\boldsymbol{\alpha}_i, \mathbf{D}) p(\boldsymbol{\alpha}_i) d\boldsymbol{\alpha}_i \right) - \log p(\mathbf{D}).
 \end{aligned} \tag{2.10}$$

The formulation is again a pointwise estimator for the dictionary \mathbf{D} but it can be argued to be more “Bayesian” due to the treatment of the latent variables $\boldsymbol{\alpha}_i$. Another possibility is to maximize the data likelihood $p(\mathbf{X}|\mathbf{D})$, resulting in the same formulation as (2.10) without the prior term $-\log p(\mathbf{D})$; this is the strategy adopted for instance by Lewicki and Olshausen [1999]. Note that optimizing (2.10) is difficult since the integrals do not admit an analytical closed form, and some approximations have to be made [Lewicki and Olshausen, 1999].

Recently, a fully Bayesian dictionary learning formulation has been proposed by Zhou et al. [2009, 2012]. By “fully Bayesian”, we mean that all model parameters such as dictionary \mathbf{D} , coefficients $\boldsymbol{\alpha}$, but also hyper-parameters such as σ are modeled with probability distributions, and the full posterior density is estimated by using Gibbs sampling. The model is in fact significantly different than the one described by (2.6) and (2.7). In a nutshell, it involves a Beta-Bernoulli process for selecting the non-zero coefficients in $\boldsymbol{\alpha}$, Gaussian priors with Gamma hyperpriors for the value of these coefficients and for the dictionary elements. Compared to the traditional matrix factorization formulation (2.2), a Bayesian treatment has both advantages and drawbacks, which often appear in passionate discussions between Bayesians and frequentists: on the one hand, the Bayesian formulation is robust to model misspecification and can learn some hyper-parameters such as the noise level σ ; on the other hand, Bayesian inference is significantly more involved and computationally costly than matrix factorization.

Convolutional sparse coding. The dictionary learning formulation of Olshausen and Field [1996, 1997] is appropriate for processing natural image patches. A natural extension to *full* images instead of patches is called “convolutional sparse coding”. It consists of linearly decomposing an image by using small dictionary elements placed at all possible locations in the image. Such a principle was described early by Zhu et al. [2005] without concrete application. It has been revisited recently, and has been shown useful for various recognition tasks [Zeiler et al., 2010, 2011, Rigamonti et al., 2011, 2013, Kavukcuoglu et al., 2010a]

Specifically, let us consider an image \mathbf{x} represented by a vector of size l and let the binary matrix \mathbf{R}_k in $\{0, 1\}^{m \times l}$ be the linear operator such that $\mathbf{R}_k \mathbf{x}$ is the patch of size $\sqrt{m} \times \sqrt{m}$ from \mathbf{x} centered at the pixel indexed by k . The purpose of \mathbf{R}_k is to “extract” a patch at a specific location, but it is also easy to show that its adjoint operator \mathbf{R}_k^\top is such that $\mathbf{R}_k^\top \mathbf{z}$, for some vector \mathbf{z} in \mathbb{R}^m , is a vector of size l representing an image with zeroes everywhere except the patch k that contains the m entries of \mathbf{z} . In other words, \mathbf{R}_k^\top positions at pixel k a small patch in the larger image. To simplify, we assume that there is no boundary effect when the pixel k is close to the image border, *e.g.*, we use zero-padding to define patches that overlap with the image boundary. Then, it becomes natural to decompose \mathbf{x} with the Lasso formulation:

$$\min_{\mathbf{A} \in \mathbb{R}^{p \times l}} \frac{1}{2} \left\| \mathbf{x} - \sum_{k=1}^l \mathbf{R}_k^\top \mathbf{D} \boldsymbol{\alpha}_k \right\|_2^2 + \lambda \sum_{k=1}^l \|\boldsymbol{\alpha}_k\|_1, \quad (2.11)$$

where \mathbf{D} is a dictionary in $\mathbb{R}^{m \times p}$, and, as usual, $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_l]$. Solving (2.11) is easy with any standard method adapted to large-scale ℓ_1 -regularized convex problems [see Bach et al., 2012a, for a review]. For instance, Rigamonti et al. [2011] uses a proximal gradient method, whereas Kavukcuoglu et al. [2010a] uses a coordinate descent scheme. To conduct the experiment of this paragraph, we have implemented the ISTA and FISTA algorithms of Beck and Teboulle [2009] for solving (2.11), and we have made it available in the SPAMS toolbox.⁶

Given now a collection of n images $\mathbf{x}_1, \dots, \mathbf{x}_n$, which are assumed

⁶<http://spams-devel.gforge.inria.fr/>.

to be of the same size for simplicity, the dictionary \mathbf{D} can be learned by minimizing

$$\min_{\substack{\{\mathbf{A}_i \in \mathbb{R}^{p \times l}\}_{i=1 \dots n} \\ \mathbf{D} \in \mathcal{C}}} \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} \left\| \mathbf{x}_i - \sum_{k=1}^l \mathbf{R}_k^\top \mathbf{D} \boldsymbol{\alpha}_{i,k} \right\|_2^2 + \lambda \sum_{k=1}^l \|\boldsymbol{\alpha}_{i,k}\|_1 \right), \quad (2.12)$$

where $\mathbf{A}_i = [\boldsymbol{\alpha}_{i,1}, \dots, \boldsymbol{\alpha}_{i,k}]$ for all $i = 1, \dots, n$. A natural optimization scheme for addressing (2.12) is to alternate between the minimization of \mathbf{D} with \mathbf{A} fixed and vice versa, as often done for the classical dictionary learning problem (see Section 5). Even though finding the global optimum of (2.12) is not feasible because the objective function is non-convex, alternate minimization provides a stationary point [Bertsekas, 1999]. Updating the dictionary \mathbf{D} with fixed coefficients \mathbf{A} can be achieved by projected gradient descent.

In Figure 2.10, we visualize two dictionaries of size $p = 100$ with elements of size $m = 16 \times 16$ pixels. We learned them on 30 whitened natural images that are rescaled such that each $\sqrt{m} \times \sqrt{m}$ patch has unit ℓ_2 -norm on average. We perform 300 steps of alternate minimization. For updating \mathbf{A} or \mathbf{D} , we always use 5 iterations of the algorithm ISTA [Beck and Teboulle, 2009], yielding a small decrease of the objective function at each step. Interestingly, we observe that the learned features have the following properties: (i) in general, they are well centered; (ii) some of them are more complex than the ones obtained with the classical dictionary learning formulation. In other words, they go beyond simple Gabor features, with curvy patterns, corners, and blobs.

2.7 Discussion

In this section, we have reviewed a large number of unsupervised learning techniques that are able to discover underlying structures in natural images. We have focused on matrix factorization and sparse coding, but other techniques that are out of the scope of our study are known to achieve similar results. We briefly mention a few of them.

A successful approach is the Gaussian mixture model (GMM). Even though it is classical, it has been shown only recently to be well

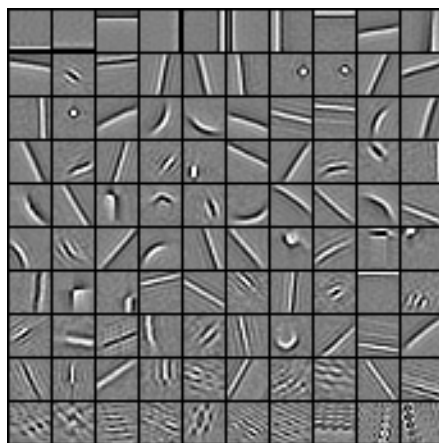
(a) With $\lambda = 0.2$.

Figure 2.10: Visualization of $p = 100$ dictionary elements learned on 30 whitened natural images. Dictionary elements are ordered from most to least used (from left to right, then top to bottom).

adapted to represent natural image patches. GMMs for image patches were first investigated by Yu et al. [2012] and then further developed by Zoran and Weiss [2011], yielding very good results for several image restoration tasks. The probabilistic model consists of representing the distribution of a patch \mathbf{x} as a convex combination of p Gaussian distributions with means $\boldsymbol{\mu}_k$ and (often diagonal) covariance matrices $\boldsymbol{\Sigma}_k$ for $k = 1, \dots, p$. Given a database of patches $\mathbf{x}_1, \dots, \mathbf{x}_n$, the EM-algorithm [Dempster et al., 1977] is typically used to learn means and covariances. When displaying the means, small localized Gabor filters typically appear when the data is whitened [Coates et al., 2011].

Finally, small oriented localized filters can also be learned from natural image patches by using undirected graphical models such as restricted Boltzmann machines (RBM) [see Hinton, 2002, Bengio, 2009, Coates et al., 2011], and also by using convolutional neural networks [LeCun et al., 1998a, Zeiler and Fergus, 2014], which have recently gained some popularity for solving large-scale visual recognition tasks.

3

Sparse Models for Image Processing

We have previously shown that natural image patches could be modeled by using dictionary learning techniques and its variants, but we have not presented any concrete application yet. We have indeed focused so far on the *interpretation* of the visual patterns obtained by various methods, but one may wonder whether or not these interpretable structures can be useful for prediction tasks. For quite a long time, the answer to this question was uncertain, until the work of Elad and Aharon [2006] on image denoising that achieved state-of-the-art results compared to other approaches at that time.

The current section is devoted to several applications of dictionary learning in image processing. First, we present the simple and yet effective scheme for image denoising introduced by Elad and Aharon [2006], and then move to more complex tasks. In general, we show that dictionary learning can be useful for predicting missing visual information, leading to state-of-the-art results for image inpainting and demosaicking [Mairal et al., 2008c, 2009], super-resolution and deblurring [Yang et al., 2010a, 2012a, Couzinie-Devy et al., 2011, Dong et al., 2011b, Zeyde et al., 2012, Wang et al., 2012], face compression [Bryt and Elad, 2008], or for inverting non-linear local trans-

formations [Mairal et al., 2012]. We also present natural extensions of dictionary learning to video processing [Protter and Elad, 2009], and finally, we conclude this section with a presentation of other patch-modeling approaches [Buades et al., 2005, Awate and Whitaker, 2006, Dabov et al., 2007a, Takeda et al., 2007, Zoran and Weiss, 2011, Yu et al., 2012, Chatterjee and Milanfar, 2012].

A small part of the material of this section is borrowed from the PhD thesis of the first author [Mairal, 2010].

3.1 Image denoising

Let us consider the classical problem consisting of restoring a noisy image \mathbf{y} in \mathbb{R}^n that has been corrupted by white Gaussian noise with known standard deviation σ .¹ In many cases, image denoising is formulated with an *energy* to minimize [*e.g.*, Rudin et al., 1992]:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 + \lambda \psi(\mathbf{x}), \quad (3.1)$$

where $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$ is a regularization function, and the quadratic “data-fitting” term ensures that the estimate \mathbf{x} is close to the noisy observation \mathbf{y} . When (3.1) is derived from a probabilistic image model, it is often interpreted from the point of view of maximum a posteriori estimation, where $\lambda \psi(\mathbf{x})$ is related to some negative log prior distribution on \mathbf{x} . Even though (3.1) looks simple at first sight, finding a good regularization function ψ is difficult, and in fact, it is probably one of the most important research topic in image processing nowadays. In early work, various smoothness assumptions about \mathbf{x} have led to different functions ψ . Specifically, natural images were assumed to have small total variation [Rudin et al., 1992], or the smoothness between adjacent pixel values was modeled with Markov random fields (MRF) [Zhu and Mumford, 1997].

As described in Section 1.2, sparse image models based on wavelets have also been popular for the denoising task [see Mallat, 2008]. This line of work has inspired the method we present in this section but

¹The formulations presented in this section can be also modified to handle other types of noise, such as Poisson noise [see Giryes and Elad, 2014].

the latter differs in two aspects: (i) the approach of Elad and Aharon [2006] is patch-based like other “modern” image processing methods [Buades et al., 2005, Dabov et al., 2007a, Roth and Black, 2009]; (ii) it adapts to the image patches with dictionary learning, and thus it does not use any pre-defined wavelet basis.

Patch denoising given a fixed dictionary. Elad and Aharon [2006] have proposed a simple denoising procedure that treats every patch independently. Let us consider the $\sqrt{m} \times \sqrt{m}$ patch \mathbf{y}_i of the noisy image \mathbf{y} , centered at the pixel indexed by i , and assume that an appropriate dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$ is given.² Then, the patch \mathbf{y}_i is denoised by the following steps:

1. center \mathbf{y}_i (see Section 2.1),

$$\mathbf{y}_i^c \triangleq \mathbf{y}_i - \mu_i \mathbf{1}_m \quad \text{with} \quad \mu_i \triangleq \frac{1}{n} \mathbf{1}_m^\top \mathbf{y}_i;$$

2. find a sparse linear combination of dictionary elements that approximates \mathbf{y}_i^c up to the noise level:

$$\min_{\boldsymbol{\alpha}_i \in \mathbb{R}^p} \|\boldsymbol{\alpha}_i\|_0 \quad \text{s.t.} \quad \|\mathbf{y}_i^c - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \leq \varepsilon, \quad (3.2)$$

where ε is proportional to the noise variance σ^2 ;

3. add back the mean component to obtain the clean estimate $\hat{\mathbf{x}}_i$:

$$\hat{\mathbf{x}}_i \triangleq \mathbf{D}\boldsymbol{\alpha}_i^* + \mu_i \mathbf{1}_m,$$

where $\boldsymbol{\alpha}_i^*$ is the solution, or approximate solution, obtained when addressing (3.2).

The approach we have just described yields two difficulties. First, it is important to choose a good value for ε . Second, problem (3.2) is unfortunately NP-hard. For patches of size $m = 8 \times 8 = 64$ pixels, Elad and Aharon [2006] recommend the value $\varepsilon = m(1.15\sigma)^2$. Another effective heuristic proposed by Mairal et al. [2009] assumes that the

²For simplicity, we always assume in this monograph that the patches are square, but all approaches can be easily extended to deal with other patch shapes.

targeted residual $\mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i$ behaves as the Gaussian noise of (supposedly known) variance σ^2 , such that the quantity $\|\mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2/\sigma^2$ follows a χ -square distribution with m degrees of freedom. Then, the heuristic sets $\varepsilon = \sigma^2 F_m^{-1}(\tau)$ where F_m^{-1} is the inverse cumulative distribution function of the χ_m^2 distribution. Selecting the value $\tau = 0.9$ leads to appropriate values of ε in practice for different patch sizes m . For dealing with (3.2), greedy algorithms such as orthogonal matching pursuit [Pati et al., 1993] are known to provide approximate solutions that are good enough for the denoising task [Elad and Aharon, 2006]. Such algorithms are presented in details in Section 5.1.

After having presented a simple “patch” denoising procedure given a fixed dictionary \mathbf{D} , we study three questions that will explain the different steps needed for denoising a full image:

1. *how do we reconstruct the full image from the local patch models?*
2. *which dictionary should we choose?*
3. *how does ℓ_1 compares with ℓ_0 for image denoising?*

From patches to full image estimation. The denoising scheme proposed by Elad and Aharon [2006] processes *independently* every patch \mathbf{y}_i from the noisy image \mathbf{y} , before constructing the denoised image by *averaging* the patch estimates $\hat{\mathbf{x}}_i$. More precisely, since the patches \mathbf{y}_i overlap, each pixel belongs to m different patches and thus admits m estimates.³ Then, the denoised image $\hat{\mathbf{x}}$ in \mathbb{R}^n is obtained as follows:

$$\hat{\mathbf{x}} = \frac{1}{m} \sum_{i=1}^n \mathbf{R}_i^\top \hat{\mathbf{x}}_i, \quad (3.3)$$

where \mathbf{R}_i in $\mathbb{R}^{n \times n}$ is a binary matrix defined as in Section 2.6—that is, $\mathbf{R}_i^\top \hat{\mathbf{x}}_i$ is a vector of size n corresponding to an image with zeroes everywhere except for the patch indexed by i that contains $\hat{\mathbf{x}}_i$. In other words, \mathbf{R}_i^\top is a linear operator that “positions” a patch of size m at

³For simplicity, we choose here to neglect boundary effects for pixels that are close to the image border, which technically belong to fewer patches.

the pixel i in a larger image with $n > m$ pixels.⁴ As a result, Eq. (3.3) is nothing else than an averaging procedure, which is the simplest way of aggregating estimators of the same quantity.

We remark that this averaging strategy is related to the translation-invariant wavelet denoising technique of Coifman and Donoho [1995], where a clean image is reconstructed by averaging estimates obtained by denoising several shifted versions of an input image. Even though aggregating estimators by straight averaging might look suboptimal, we are not aware of any other technique, in the context of local sparse linear models, leading to better results for reconstructing the final image from the estimation of overlapping patches.

Use of generic versus global versus image-adaptive dictionary. The first baseline considered by Elad and Aharon [2006] uses a fixed over-complete discrete cosine transform (DCT) dictionary, which was presented earlier in Figure 2.4(a). We call such an approach “generic” since the dictionary is pre-defined. The second strategy, dubbed “global”, consists of learning a dictionary on an external database of clean patches, simply following the methodology of Section 2. Finally, we also present the “adaptive” strategy of Elad and Aharon [2006], which learns the dictionary from the pool of noisy patches:

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \psi(\boldsymbol{\alpha}_i) \quad \text{s.t.} \quad \|\mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \leq \varepsilon, \quad (3.4)$$

where \mathcal{C} is defined in Section 2, and ψ is a sparsity-inducing penalty. Elad and Aharon [2006] originally use the ℓ_0 -penalty in place of ψ , and optimize (3.4) with the K-SVD algorithm [Aharon et al., 2006]. Later, additional experiments have shown that learning the dictionary with ℓ_1 could bring some benefits [Mairal et al., 2009]. In fact, we experimentally demonstrate in the sequel that ℓ_1 consistently provides better results for the denoising task, *when used for learning the dictionary only*.

⁴Note that the original averaging scheme of Elad and Aharon [2006] slightly differs from (3.3) since their estimate $\hat{\mathbf{x}}$ is defined as $\hat{\mathbf{x}} = (1-\lambda)\mathbf{y} + \lambda(1/m) \sum_{i=1}^n \mathbf{R}_i^\top \hat{\mathbf{x}}_i$, where \mathbf{y} is the noisy image. We have empirically found that the setting $\lambda = 1$ leads in fact to very good results and we have thus omitted the parameter λ .

Indeed, we observe that regardless of the way the dictionary is learned, ℓ_0 should always be used for the final reconstruction step in (3.2).

This conclusion is drawn from the following experiment. We consider 12 classical images presented in Figure 3.1, which were used in other benchmarks about image denoising [*e.g.*, Mairal et al., 2009]. We corrupt every image with white Gaussian noise of standard deviation σ in $\{5, 10, 15, 20, 25, 50, 100\}$ for pixel values in the range $[0; 255]$, and we measure the reconstruction quality obtained by the three approaches by using the PSNR criterion.⁵ We consider dictionaries of size $p = 256$ elements, different patch sizes $m = l \times l$, with l in $\{6, 8, 10, 12, 14, 16\}$, and we use the parameter $\tau = 0.9$ suggested before for choosing the reconstruction threshold ε . For every noise level, the parameter l is selected such that it maximizes the average PSNR obtained on the last 5 images of the dataset.

We make a comparison between the six scenarios described below. In all cases, the dictionaries are learned with the SPAMS software, by performing 10 passes over the available training data.

1. DCT: we use an overcomplete DCT dictionary; the patches are reconstructed by using the ℓ_0 -regularization (3.2).
2. ℓ_0 -global/ ℓ_0 : the dictionary is learned on 400 000 natural image patches extracted from the Kodak PhotoCD images⁶. Denoting by $\mathbf{x}_1, \dots, \mathbf{x}_n$ these patches, we learn \mathbf{D} by minimizing

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}_i\|_0 \leq s,$$

where $s = 10$ is known to approximate well clean natural image patches [Elad and Aharon, 2006].

3. ℓ_0 -adapt/ ℓ_0 : the dictionary is learned by minimizing (3.4) with $\psi = \ell_0$; we initialize the learning procedure with the global dictionary obtained in the scenario ℓ_0 -global/ ℓ_0 .

⁵Denoting by MSE the mean squared-error for images whose intensities are between 0 and 255, the PSNR is defined as $\text{PSNR} = 10 \log_{10}(255^2/\text{MSE})$ and is measured in dB. A gain of 1 dB reduces the MSE by approximately 20%.

⁶The dataset is available here: <http://r0k.us/graphics/kodak/>.

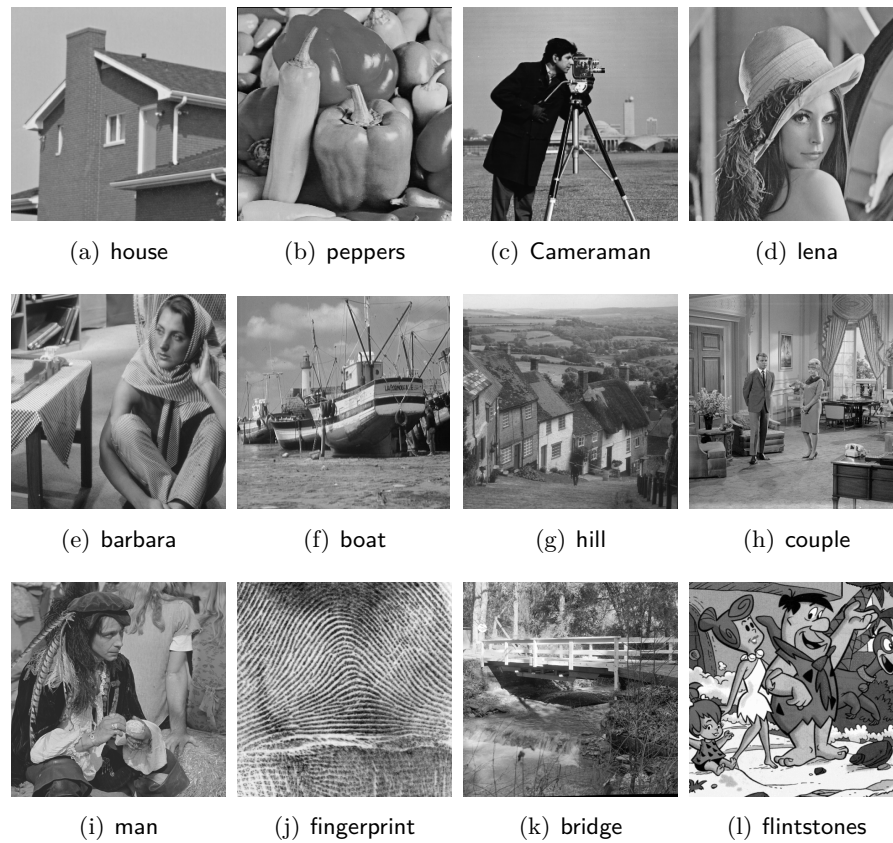


Figure 3.1: Dataset of 12 standard images used in the image denoising benchmarks.

4. ℓ_1 -global/ ℓ_0 : the scenario is the same as ℓ_0 -global/ ℓ_0 , except that the dictionary is learned with the ℓ_1 -penalty, following exactly the methodology of Section 2.4. The final reconstruction is obtained with ℓ_0 as in (3.2).
5. ℓ_1 -adapt/ ℓ_0 : same as ℓ_1 -adapt/ ℓ_0 , except that the dictionary is learned by minimizing (3.4) with $\psi = \ell_1$; we use the dictionary obtained in the scenario ℓ_1 -global for initializing the learning procedure. The final reconstruction is still obtained with the ℓ_0 -penalty.
6. ℓ_1 -adapt/ ℓ_1 : same as ℓ_1 -adapt/ ℓ_0 , but the patches \mathbf{y}_i^c in the final reconstruction are denoised by replacing the ℓ_0 -penalty by the ℓ_1 -norm in (3.2).

We report the mean PSNR obtained by the different scenarios in Table 3.1, along with the performance achieved by other state-of-the-art approaches of the literature. The conclusions are the following:

- the simple denoising scheme that we have presented performs slightly worse than more recent methods such as Dabov et al. [2007a], Mairal et al. [2009], Chatterjee and Milanfar [2012] in terms of PSNR;
- adaptive dictionaries yield better results than global ones, as already observed by Elad and Aharon [2006];
- learning the dictionary with ℓ_1 consistently yields better results than ℓ_0 , even though the ℓ_0 -penalty is used for the final image reconstruction step in (3.2);
- ℓ_0 yields better results than ℓ_1 in the final image reconstruction step, regardless of the way the dictionary was learned.

Even though the conclusion that one should “first learn the dictionary with ℓ_1 , before using ℓ_0 for denoising the image” might seem counterintuitive, we believe that the reason of the good performance of ℓ_1 for dictionary learning might be a better stability of the sparsity patterns than the ones obtained with ℓ_0 . This may subsequently

	σ						
	5	10	15	20	25	50	100
DCT	37.30	33.38	31.24	29.75	28.66	25.24	22.00
ℓ_0 -global/ ℓ_0	37.21	33.37	31.41	30.01	28.95	25.65	22.44
ℓ_1 -global/ ℓ_0	37.22	33.50	31.56	30.17	29.13	25.77	22.53
ℓ_0 -adapt/ ℓ_0	37.49	33.75	31.70	30.40	29.33	26.04	22.64
ℓ_1 -adapt/ ℓ_0	37.60	33.90	31.90	30.51	29.43	26.20	22.72
ℓ_1 -adapt/ ℓ_1	36.58	32.85	30.77	29.29	28.14	24.47	21.83
GSM	37.05	33.34	31.31	29.91	28.84	25.66	22.80
K-SVD	37.42	33.62	31.58	30.18	29.10	25.61	22.10
BM3D	37.62	34.00	32.05	30.73	29.72	26.38	23.25
LSSC	37.67	34.06	32.12	30.78	29.74	26.57	23.39
Plow	37.38	32.98	31.38	30.13	29.30	26.38	23.24
EPLL	37.36	33.64	31.67	30.32	29.29	26.12	23.03
CSR	37.61	34.00	32.05	30.72	29.70	26.53	23.45
BM3D-PCA	37.79	34.20	32.27	30.94	29.92	26.75	23.16

Table 3.1: Denoising performance in PSNR for various methods on the 12 standard images of Figure 3.1 for various levels of noise σ . GSM refers to the Gaussian scale mixture model of Portilla et al. [2003]; K-SVD refers to the original method of Elad and Aharon [2006]; BM3D refers to the state-of-the-art denoising approach of Dabov et al. [2007a]; LSSC refers to Mairal et al. [2009], Plow refers to Chatterjee and Milanfar [2012], EPLL to Zoran and Weiss [2011], CSR to Dong et al. [2013], and BM3D-SAPCA to an extension of BM3D with better results [Dabov et al., 2009, Katkovnik et al., 2010]. For all these approaches, publicly available software is available on the corresponding authors' web pages.

yield a better behavior in terms of optimization. Whereas the ℓ_1 -scheme guarantees us to obtain a stationary point of the dictionary learning formulation, the ℓ_0 -counterpart does not. A similar experiment was conducted by Mairal [2010, chapter 1.6], with similar conclusions.

3.2 Image inpainting

Dictionary learning is well adapted to the presence of missing data—that is, it is appropriate for *inpainting* [Bertalmio et al., 2000] when the missing pixels form small holes that are smaller than the patch

sizes. To deal with unobserved information, the dictionary learning formulation can be modified by introducing a binary mask \mathbf{M}_i for every patch indexed by i [Mairal et al., 2008c,e]. Formally, we define \mathbf{M}_i as a diagonal matrix in $\mathbb{R}^{m \times m}$ whose value on the j -th entry of the diagonal is 1 if the pixel $\mathbf{y}_i[j]$ is observed and 0 otherwise. Then, the dictionary learning formulation becomes

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{M}_i(\mathbf{y}_i - \mathbf{D}\boldsymbol{\alpha}_i)\|_2^2 + \lambda\psi(\boldsymbol{\alpha}_i),$$

where ψ is a sparsity-inducing penalty, $\mathbf{M}_i\mathbf{y}_i$ represents the observed pixels from the i -th patch of the image \mathbf{y} and $\mathbf{D}\boldsymbol{\alpha}_i$ is the estimate of the full patch i . In practice, the binary mask does not drastically change the optimization procedure, and one can still use classical optimization techniques for dictionary learning, *e.g.*, alternating between the optimization of \mathbf{D} and \mathbf{A} , as described in Section 5. When the image \mathbf{y} is only corrupted by missing pixels and not by other noise source, it is also possible to enforce hard reconstruction constraints, leading to the formulation

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \sum_{i=1}^n \psi(\boldsymbol{\alpha}_i) \quad \text{s.t.} \quad \mathbf{M}_i\mathbf{y}_i = \mathbf{M}_i\mathbf{D}\boldsymbol{\alpha}_i. \quad (3.5)$$

Similarly to the denoising problem, once the dictionary is learned, we need to sparsely encode all overlapping patches, *e.g.*, by approximately minimizing (3.5) with a greedy algorithm when ψ is the ℓ_0 -penalty. Then, the full image is obtained by averaging using (3.3).

Before showing any inpainting result, we shall comment on *when the formulation described here is supposed to work*. Our first remark is that it can only handle holes that are smaller than the patch size. Dealing with larger holes might be possible, but with a different formulation that would be able to synthesize new information in empty patches [see, *e.g.* Criminisi et al., 2004, Peyré, 2009, Roth and Black, 2009]. Second, one assumes that *the noise pattern is unstructured*, such that no noise patterns are learned by the dictionary. The demosaicking task from the next section is a typical example with very structured missing patterns, and where an alternative strategy needs to be used.

Finally, we show inpainting results in Figure 3.2. For both images, the noise structure is relatively random when seen at the patch level, and the dictionary learning approach performs well at recovering the missing pixels. In particular, the brick texture in the image *house* is hardly visible for the human eye, but it is well recovered by the algorithm.

3.3 Image demosaicking

The problem of demosaicking consists of reconstructing a color image from the raw information provided by colored-filtered sensors from digital cameras [Kimmel, 1999]. Most consumer-grade cameras use a grid of sensors; for every pixel, one sensor records the amount of light it receives for a short period of time. A color filter, red, green, or blue, is placed in front of the sensor such that only one color channel is measured for each pixel. Reconstructing the missing channels for the full image is a task called *demosaicking*. It is usually performed on-the-fly by digital cameras, but can also be achieved with a better quality by professional software on a desktop computer.

A typical pattern of filters, called the “Bayer” pattern, is displayed in Figure 3.3. Every row alternates between red and green, or green and blue filters. As a result, information captured by digital sensors in the green channel is twice more important than in the red or blue channels, even though the final post-processed color image observed by the user does not reflect this fact. An example of a mosaicked image with the Bayer pattern is presented in Figure 3.4(a). Reconstructing the color image without producing any visual artifact is difficult since it requires recovering frequencies in every channel that are greater than the sampling rate. Specifically, the resolution in the red and blue channels needs to be doubled.

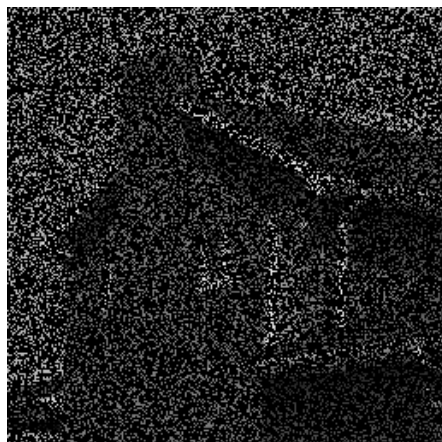
Since demosaicking consists of recovering missing information for every pixel, it can be cast as an inpainting problem, but the inpainting procedure described in the previous section cannot be directly applied because the noise pattern is strongly structured. An effective two-step strategy proposed by Mairal et al. [2008c] to overcome the shortcoming



(a) Example A, Damaged



(b) Example A, Restored



(c) Example B, Damaged



(d) Example B, Restored

Figure 3.2: Top: inpainting result obtained by using the source code of Mairal et al. [2008c], where the text is automatically removed on the restored image. Bottom: inpainting result from Mairal et al. [2008e], where 80% of the pixels are randomly missing from the original image. The algorithm is able to reconstruct the brick texture on the right, by adapting the dictionary to the available information displayed on the left. The bottom images are under “copyright ©2008 Society for Industrial and Applied Mathematics. Reprinted with permission. All rights reserved”. Best seen by zooming on a computer screen.

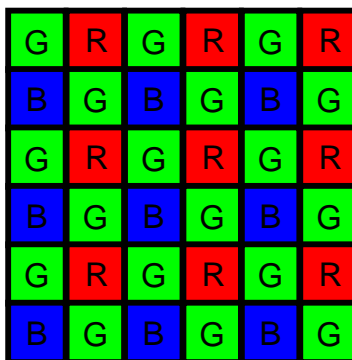


Figure 3.3: Common Bayer pattern used for digital camera sensors.

of the original inpainting formulation consists of the following steps:

1. learn a dictionary \mathbf{D}_0 on an external database of color patches;
2. use \mathbf{D}_0 to obtain a first estimate of the color image, following the image reconstruction procedure of the previous section;
3. retrain the dictionary on the first image estimate, leading to a new dictionary \mathbf{D}_1 adapted to the image content.
4. use \mathbf{D}_1 to obtain the final estimate of the color image as in 2.

In addition, the patches should also be centered according to the recommendations of Section 2.1. The results achieved by this method can be were shown by Mairal et al. [2008c] to outperform the state of the art in terms of PSNR on the 24 images from the Kodak PhotoCD benchmark, even though some visual artifacts remain visible in areas that are particularly difficult to reconstruct, as shown in Figure 3.4(b).

The demosaicking approach based on dictionary learning has been improved later by Mairal et al. [2009] by combining sparse estimation with image self-similarities [Buades et al., 2009]. The resulting method yields a higher PSNR, but more importantly it seems to significantly reduce visual artifacts (see Figure 3.4(c)), and achieves state-of-the-art results in terms of PSNR [Menon and Calvagno, 2011].

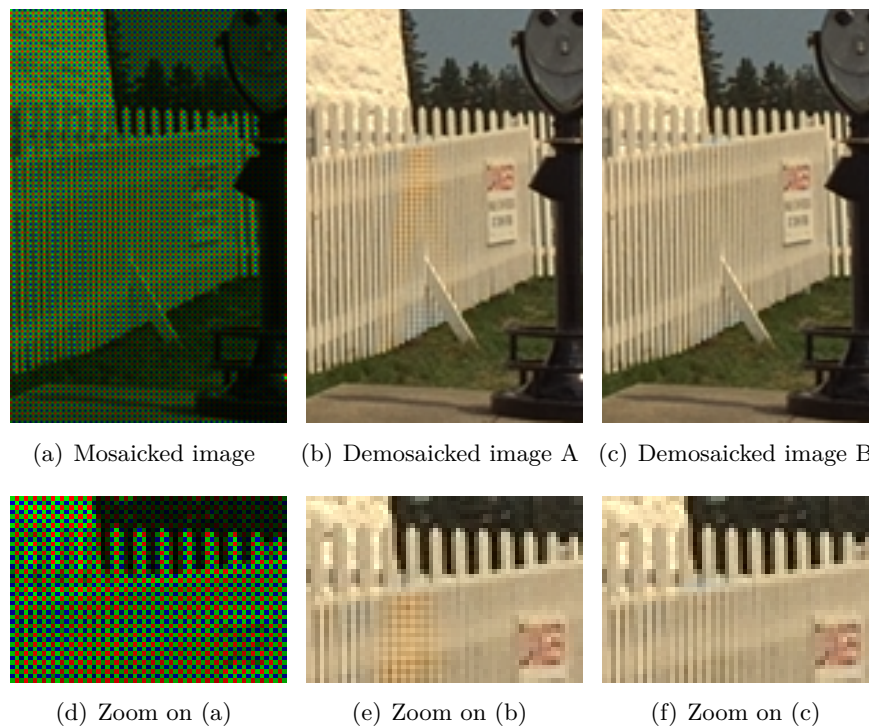


Figure 3.4: (a): Mosaicked image; (b): image demosaicked with the approach of Mairal et al. [2008c] presented in Section 3.3. (c): image demosaicked with the non-local sparse image models of Mairal et al. [2009]. Note that the part of the image displayed here is known to be particularly difficult to reconstruct. Other reconstructions obtained from the database of 24 images is artifact-free in general. Best seen in color by zooming on a computer screen.

3.4 Image up-scaling

Dictionary learning has also gained important success for reconstructing high-resolution images from low-resolution ones, a problem called *image up-scaling* or *digital zooming*.⁷ The technique originally proposed by Yang et al. [2010a] consists of learning a pair of dictionaries $(\mathbf{D}_l, \mathbf{D}_h)$ that can respectively reconstruct low- and high-resolution patches with the same sparse decomposition coefficients. Then, the relation between \mathbf{D}_l and \mathbf{D}_h is exploited for processing new low-resolution images and turn them into high-resolution ones.

Similar to the “example-based” approach of Freeman et al. [2002], the method requires a database of pairs of training patches $(\mathbf{x}_i^l, \mathbf{x}_i^h)_{i=1}^n$, where \mathbf{x}_i^l in \mathbb{R}^{m_l} is a low-resolution version of the patch \mathbf{x}_i^h in \mathbb{R}^{m_h} , and $m_h > m_l$. The database is built from training images that are generic enough to represent well the variety of natural images. Then, the pairs of patches are used offline to train the dictionaries \mathbf{D}_l in $\mathbb{R}^{m_l \times p}$ and \mathbf{D}_h in $\mathbb{R}^{m_h \times p}$, which are subsequently used for image up-scaling in a way that will be described in the sequel. The methodology is related to the *global* approach for image denoising from Section 3.1, where the dictionary is not adapted to the image at hand that needs to be processed, but is adapted instead to a generic database of images.

The image upscaling method of Yang et al. [2010a] has received a significant amount of attention because of the quality of its results, and has been improved recently in various ways [Couzinie-Devy et al., 2011, Zeyde et al., 2012, Dong et al., 2011b, Wang et al., 2012, Yang et al., 2012a]. In this section, we briefly review some of the ideas that have appeared in this line of work, but refer to the corresponding papers for more details. We start with the original approach of Yang et al. [2010a], before moving to two natural variants.

The original approach of Yang et al. [2010a]. The dictionary learning formulation of Yang et al. [2010a] for image up-scaling consists of

⁷The terminology of super-resolution is sometimes used but leads to some confusion. “Super-resolution” traditionally means synthesizing high-resolution images from a sequence of low-resolution ones [see the discussion in Elad, 2010, page 341]. Here, we are dealing with a single input image.

learning \mathbf{D}_l and \mathbf{D}_h with the following joint minimization problem:

$$\min_{\substack{\mathbf{D}_l \in \mathcal{C}_l \\ \mathbf{D}_h \in \mathbb{R}^{m_h \times p} \\ \mathbf{A} \in \mathbb{R}^{p \times n}}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2m_l} \left\| \mathbf{x}_i^l - \mathbf{D}_l \boldsymbol{\alpha}_i \right\|_2^2 + \frac{1}{2m_h} \left\| \mathbf{x}_i^h - \mathbf{D}_h \boldsymbol{\alpha}_i \right\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1, \quad (3.6)$$

where $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n]$ is the matrix of sparse coefficients, and \mathcal{C}_l is the set of matrices in $\mathbb{R}^{m_l \times p}$ whose columns have less than unit ℓ_2 -norm. The formulation encourages the dictionaries to provide a common sparse representation $\boldsymbol{\alpha}_i$ for the pair of patches \mathbf{x}_i^l and \mathbf{x}_i^h . Note that, to simplify, we have omitted pre-processing steps for the patches \mathbf{x}_i^l , which lead to substantially improved results in practice [see Yang et al., 2010a, Zeyde et al., 2012], and we have preferred to focus on the main principles of the method.

Once the dictionaries are trained, they can be used for turning new low-resolution images into high-resolution ones. Specifically, let us consider an image \mathbf{y}^l of size $\sqrt{n_l} \times \sqrt{n_l}$, which is independent from the training database, and assume that the goal is to synthesize a high-resolution image \mathbf{y}^h of larger size $\sqrt{n_h} \times \sqrt{n_h}$ —say, for instance, $\sqrt{n_h} = 2\sqrt{n_l}$ when one wishes to double the resolution.

According to (3.6), the main underlying assumption is that high-resolution and low-resolution variants of the same patch admit a common sparse decomposition onto \mathbf{D}_h and \mathbf{D}_l , respectively. Then, given a patch \mathbf{y}_i^l in \mathbb{R}^{m_l} centered at pixel i in the low-resolution image \mathbf{y}^l , the method of Yang et al. [2010a] computes a sparse vector $\boldsymbol{\beta}_i$ in \mathbb{R}^p such that $\mathbf{y}_i^l \approx \mathbf{D}_l \boldsymbol{\beta}_i$, and the high-resolution estimate of the patch is simply $\mathbf{D}_h \boldsymbol{\beta}_i$. Such an approach is however not fully satisfactory regarding the formulation (3.6), where each sparse decomposition is obtained by using both the low- and high-resolution patches. In other words, $\boldsymbol{\beta}_i$ should be ideally computed by minimizing

$$\min_{\boldsymbol{\beta}_i \in \mathbb{R}^p} \frac{1}{2m_l} \left\| \mathbf{y}_i^l - \mathbf{D}_l \boldsymbol{\beta}_i \right\|_2^2 + \frac{1}{2m_h} \left\| \mathbf{y}_i^h - \mathbf{D}_h \boldsymbol{\beta}_i \right\|_2^2 + \lambda \|\boldsymbol{\beta}_i\|_1, \quad (3.7)$$

but unfortunately only \mathbf{y}_i^l is available at test time since \mathbf{y}_i^h is the unknown quantity to estimate. Such a discrepancy is addressed by Yang et al. [2010a] with a heuristic, where \mathbf{y}_i^h in (3.7) is replaced by

an auxiliary variable representing the current prediction of the high-resolution patches given previously processed neighbors. Finally, reconstructing the full image from the local patch estimates can be done as usual by straight averaging, even though other possibilities yielding good results have also been proposed [see Zeyde et al., 2012].

First variation with a regression formulation. A strategy for overcoming the discrepancy between training and testing in the previous approach has been independently proposed by several authors [Zeyde et al., 2012, Couzinie-Devy et al., 2011, Yang et al., 2012a]. In all these approaches, the sparse coefficients α_i and β_i are always computed by using *low-resolution patches only*. This motivates the following two-step training approach:

1. compute \mathbf{D}_l and \mathbf{A} with a classical dictionary learning formulation, for instance:

$$\min_{\mathbf{D}_l \in \mathcal{C}_l, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i^l - \mathbf{D}_l \alpha_i\|_2^2 + \lambda \|\alpha_i\|_1.$$

2. obtain \mathbf{D}_h by solving a multivariate *regression* problem:

$$\min_{\mathbf{D}_h \in \mathbb{R}^{m_h \times p}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i^h - \mathbf{D}_h \alpha_i\|_2^2,$$

where the α_i 's are fixed after the first step. See also Zeyde et al. [2012] for other variants.

Once the dictionaries are learned, a new low-resolution image \mathbf{y}_l is processed by decomposing its patches \mathbf{y}_i^l onto \mathbf{D}_l , yielding sparse coefficients β_i , and the estimate of the corresponding high resolution patches are again the quantities $\mathbf{D}_h \beta_i$.

Second variation with bilevel optimization. The previous variant addresses one shortcoming of the original approach of Yang et al. [2010a], but unfortunately loses the ability of jointly learning \mathbf{D}_l and \mathbf{D}_h . The variant presented in this paragraph is motivated by that issue. It was independently proposed by Couzinie-Devy et al. [2011] and Yang et al.

[2012a] and allows learning \mathbf{D}_l and \mathbf{D}_h at the same time. As we shall see, it raises a challenging bilevel optimization problem.

Let us first introduce a notation describing the solution of the Lasso for any vector \mathbf{x} and dictionary \mathbf{D} :

$$\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D}) \triangleq \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \right],$$

where we assume the solution of the Lasso problem to be unique.⁸ Then, the joint dictionary learning formulation consists of minimizing

$$\min_{\substack{\mathbf{D}_l \in \mathcal{C}_l \\ \mathbf{D}_h \in \mathbb{R}^{m_h \times p}}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \left\| \mathbf{x}_i^h - \mathbf{D}_h \boldsymbol{\alpha}^*(\mathbf{x}_i^l, \mathbf{D}_l) \right\|_2^2. \quad (3.8)$$

The problem is challenging since it is non-convex, non-differentiable, and the dependency of the objective with respect to \mathbf{D}_l goes through the argmin of the Lasso formulation. Such a setting is sometimes referred to as “bilevel optimization” problems and is known to be difficult. Eq. (3.8) appears under the name of “coupled dictionary learning” in [Yang et al., 2012a], and is in fact a particular case of the more general “task-driven dictionary learning” formulation of Mairal et al. [2012], which we will present in more details in Section 4.5 about visual recognition. It is possible to show that in the asymptotic regime where n grows to infinity, the cost function to optimize becomes differentiable and that its gradient can be estimated. Based upon this observation, a stochastic gradient descent algorithm can be used as an effective heuristic for finding an approximate solution [see Couzinie-Devy et al., 2011, Mairal et al., 2012]. We conclude this section with a few visual results, which we present in Figure 3.5.

3.5 Inverting nonlinear local transformations

Interestingly, the image up-scaling formulations that we have presented previously do not explicitly take into account the type of transformation between the patches \mathbf{x}_i^h and \mathbf{x}_i^l . In other words, they only assume

⁸When the solution is not unique, the elastic-net regularization [Zou and Hastie, 2005] can be used instead of the ℓ_1 -norm. It consists of replacing the penalty $\|\boldsymbol{\alpha}\|_1$ by $\|\boldsymbol{\alpha}\|_1 + (\mu/2)\|\boldsymbol{\alpha}\|_2^2$. The resulting penalty still enjoy sparsity-inducing properties, but is strongly convex and thus yields a unique solution.

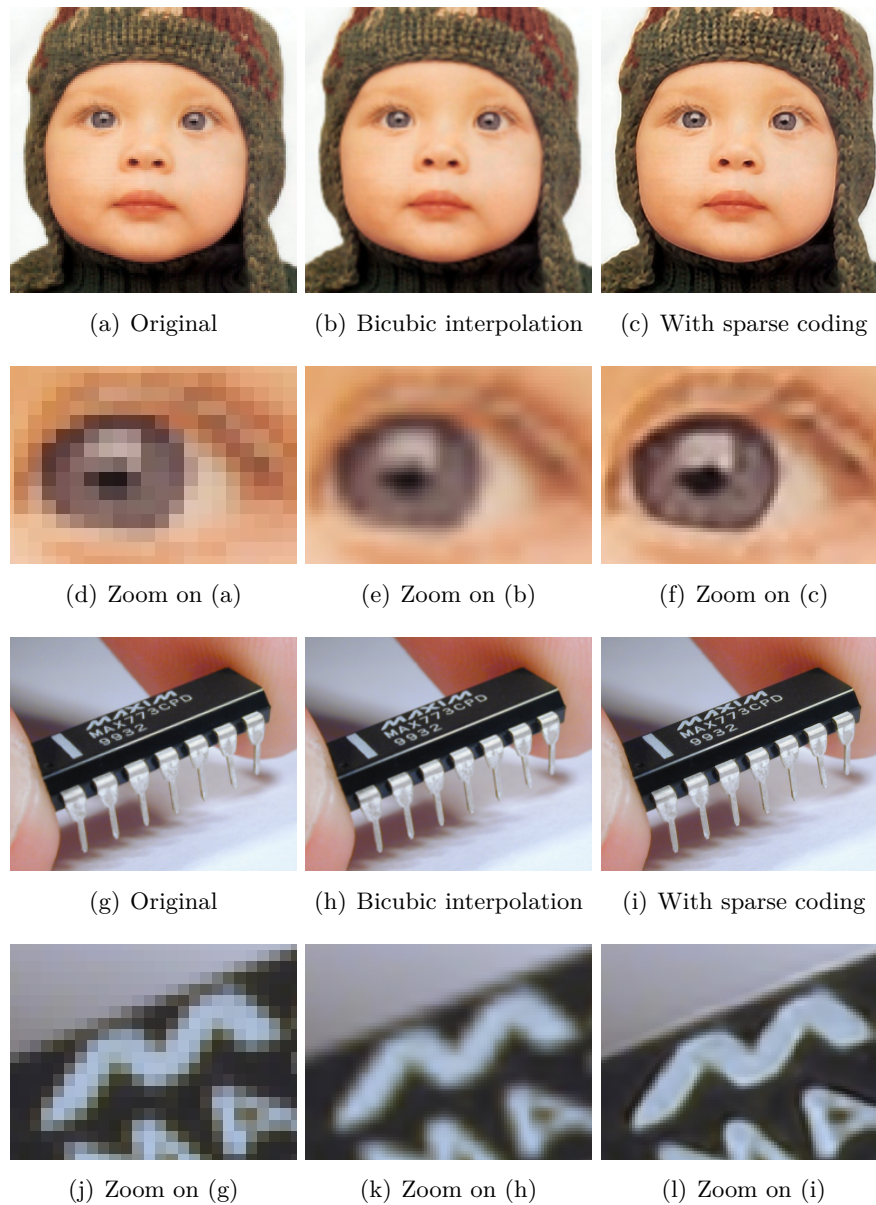


Figure 3.5: Image upscaling results obtained with the dictionary learning technique of Couzinie-Devy et al. [2011]. Images produced by Florent Couzinie-Devy.

that the pairs of patches \mathbf{x}_i^h and \mathbf{x}_i^l admit a common sparse representation onto two different dictionaries, but they do not make any a priori assumption on the relation between the patches. For this reason, it is appealing to try similar formulations for other problems than image upscaling.

Several extensions have been developed along this line of work. Couzinie-Devy et al. [2011] and Dong et al. [2011b] have for instance developed image deblurring techniques that are effective when the blur is local. Other less-standard applications also appear in the literature. Dong et al. [2011b] propose indeed to learn a mapping between face photographs and hand-drawn sketches. They use a database of images called CUFS [Wang and Tang, 2009] of about 600 subjects. For each subject, a grayscale photograph is provided, along with a sketch drawn by an artist. The dataset can be used to create pairs of photograph-sketch patches and a nonlinear mapping is learned between the two patch spaces, providing convincing visual results.

Finally, we illustrate the ability of the previous formulations for learning nonlinear mappings with the problem of reconstructing grayscale images from binary ones, a task called “inverse halftoning”. With the use of binary display and printing technologies in the 70’s, converting a grayscale continuous-tone image into a binary one that looks perceptually similar to the original one (“halftoning”) was of utmost importance. A classical algorithm was developed for that purpose by Floyd and Steinberg [1976]. The goal of “inverse halftoning” is to restore these binary images and estimate the original ones. Since building a large database of pairs of grayscale-halftoned image patches is easy, the methodology of the previous section can be used for learning a mapping between binary and grayscale patches. In Figure 3.6, we display a few visual results obtained with the approach of Mairal et al. [2012]. The dictionaries were trained on an external database build with the Floyd-Steinberg algorithm.

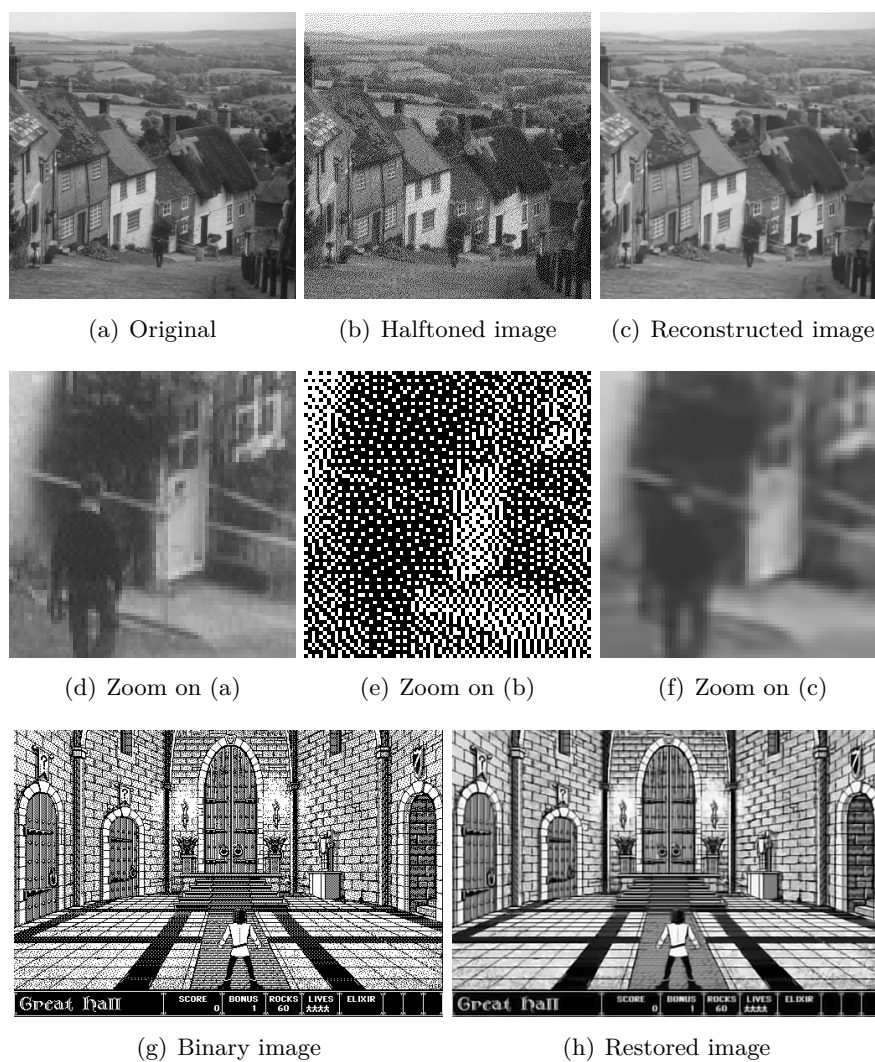


Figure 3.6: Inverse halftoning results obtained by Mairal et al. [2012]. The last row presents a result on a binary image publicly available on the Internet. No ground truth is available for this binary image from an old computer game. Despite the fact that the algorithm that has generated this image is unknown and probably different than the Floyd-Steinberg algorithm, the result is visually convincing. Best seen by zooming on a computer screen.

3.6 Video processing

Extensions of dictionary learning techniques for dealing with videos have been proposed by Protter and Elad [2009]. Given a noisy video sequence, the most naive approach consists of processing each frame independently. However, significantly better results can be obtained by exploiting the temporal consistency that exists among consecutive frames. To do so, a few modifications to the classical denoising formulation are sufficient:

- several frames can be processed at the same time by considering three-dimensional patches that involve one temporal dimension. For instance, a video patch can be of size $m = e \times e \times T$, where T is the number of frames in the patch, *e.g.*, $e = 10$ pixels and $T = 5$ frames.
- After processing T frames starting at time t , we obtain a dictionary \mathbf{D}_t adapted to the set of three-dimensional patches at time t . When moving to the next block of T frames at time $t + 1$ (considering that the blocks overlap in time), we can learn a dictionary \mathbf{D}_{t+1} adapted to the information available at time $t + 1$, using \mathbf{D}_t as an initialization to the learning algorithm.

In Figures 3.7 and 3.8, we present two video processing results from [Mairal et al., 2008e], where the original video extension of Protter and Elad [2009] has been adapted to the inpainting and color video denoising tasks. The results of the “video processing” approach are significantly better than those obtained by processing independently each frame.

3.7 Face compression

We have shown so far that dictionary learning was appropriate for many *restoration* tasks, but an intuitive use of sparsity is for *compression*. Following this insight, Bryt and Elad [2008] have proposed a technique for encoding photographs of human faces. They achieve impressive compression results, where each image is represented with less

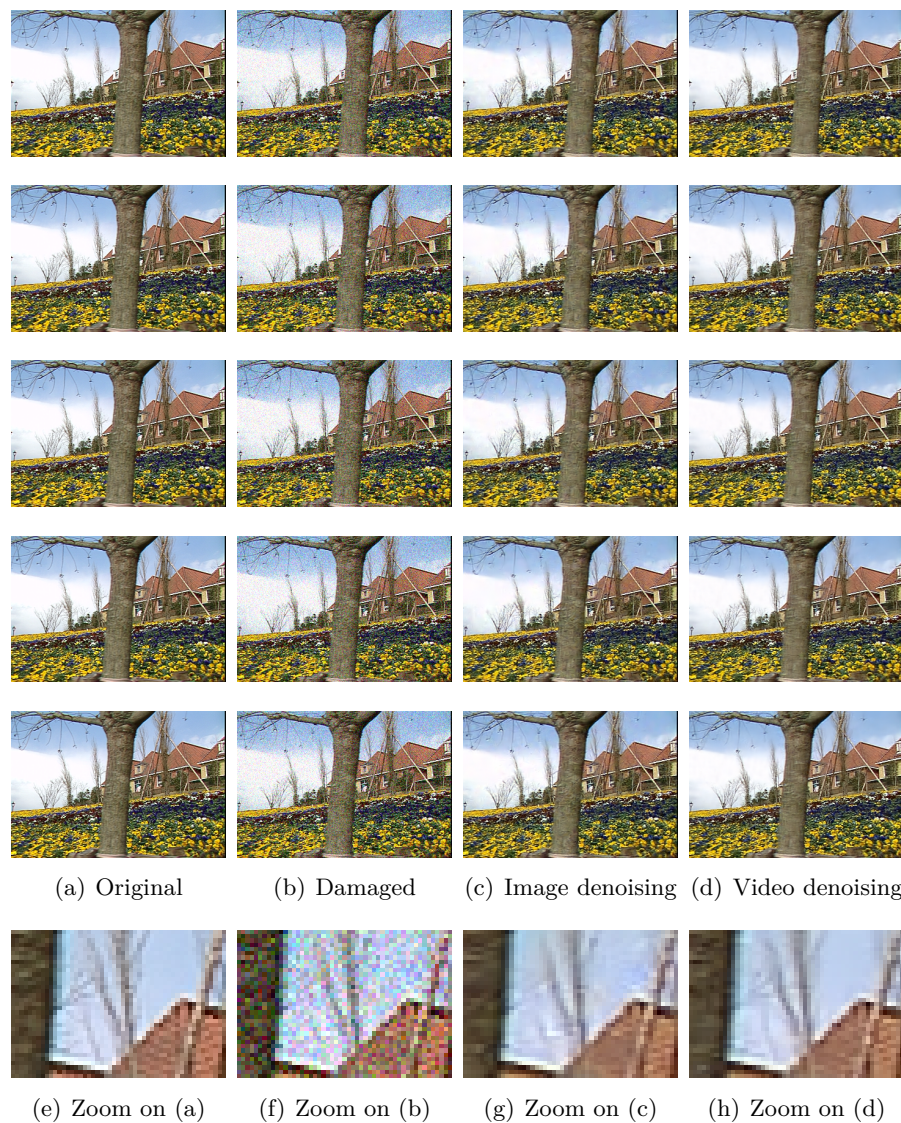


Figure 3.7: Color video denoising result from Mairal et al. [2008e]. The third column shows the result when each frame is processed independently from the others. Last column shows the result of the video processing approach. Best seen by zooming on a computer screen. “Copyright ©2008 Society for Industrial and Applied Mathematics. Reprinted with permission. All rights reserved”.

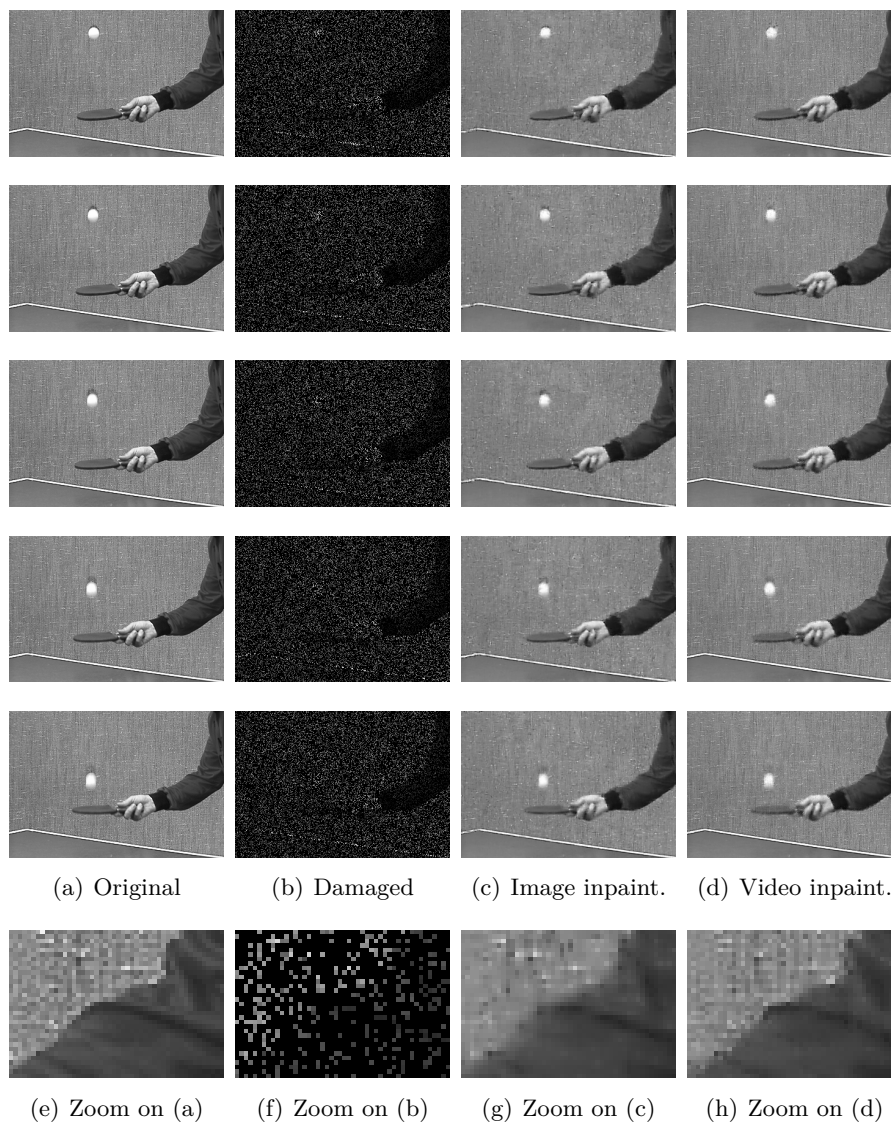


Figure 3.8: Video inpainting result from Mairal et al. [2008e]. The third column shows the result when each frame is processed independently from the others. Last column shows the result of the video processing approach. Best seen by zooming on a computer screen. “Copyright ©2008 Society for Industrial and Applied Mathematics. Reprinted with permission. All rights reserved”.

than 1000 bytes while suffering from minor loss only in visual quality. The algorithm requires a database of training face images that are processed as follows:

1. input images are geometrically aligned and downsampled to the same size. As a result, all the main face features are at the same position across images. This step is automatically performed by first detecting 13 face features and warping the images to achieve the desired alignment [see Bryt and Elad, 2008];
2. the images are sliced into 15×15 *non-overlapping* patches;
3. *for each patch position*, a dictionary of size $p = 512$ elements is learned by using all patches available at the same position from the training images. The K-SVD algorithm is chosen by Bryt and Elad [2008] for that purpose.

After this training phase, a new test image can be encoded as follows:

1. the test image goes through the same geometrical alignment process as for the training images. Encoding the inverse transformation requires 20 bytes;
2. the image is sliced into non-overlapping patches that are encoded using orthogonal matching pursuit (see Section 5.1) with a few number s of non-zero coefficients (typically $s = 4$);
3. the indices of the sparse coefficients are encoded with a Huffman table [see MacKay, 2003]; the weights are quantized into 7 bits.

Once encoded, each patch can be decoded with a simple matrix vector multiplication, and the inverse geometrical warping is applied to obtain an approximation of the original image. Visual results are presented in Figure 3.9 for face images represented by 550 bytes only. Unlike generic approaches such as JPEG and JPEG2000, the approach of Bryt and Elad [2008] can exploit the fact that each patch represents a specific information, *e.g.*, eye, mouth, which can be learned by using adaptive dictionaries. Since the patches do not overlap, the reconstructed images suffer from blocking artifacts, which can be significantly reduced by using a deblocking post-processing algorithm.

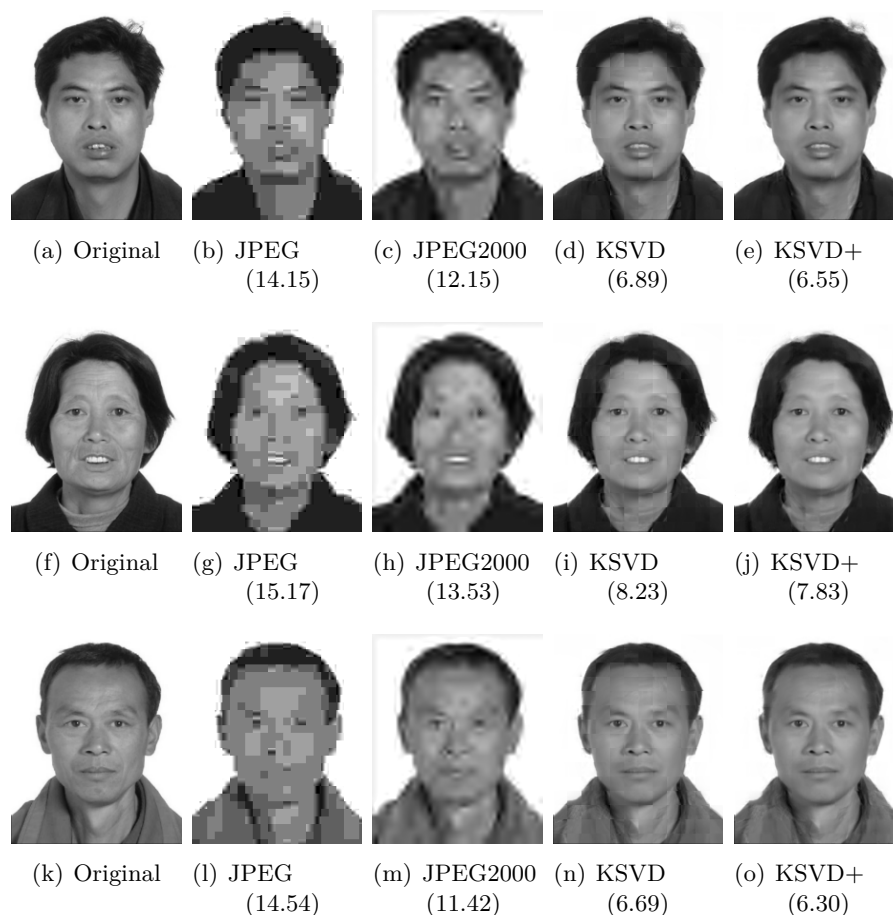


Figure 3.9: Results of the face compression method of Bryt and Elad [2008] compared to two classical image compression algorithms. All images are compressed into 550 bytes representations. The raw output of the method, denoted by KSVD, is presented in the penultimate column. The last column KSVD+ is a post-processed result obtained by processing the output of KSVD with a deblocking algorithm that reduces visual artifacts. The numbers in parenthesis are the square root of the mean squared errors (RMSE), obtained by comparing the compressed images with the original one (first column). Images kindly provided to us by Ori Bryt and Michael Elad. Best seen by zooming on a computer screen.

3.8 Other patch modeling approaches

Dictionary learning was successful in many image processing tasks because of its ability to model well natural image patches. Other patch-based methods can be used for that purpose, and have also shown impressive results for image restoration. For the sake of completeness, we briefly present a few of them that have gained some success, but we refer the reader to the corresponding papers for more details.

Non-local means and non-parametric approaches. Efros and Leung [1999] showed that self-similarities inherent to natural images could be used effectively in texture synthesis tasks. Following their insight, Buades et al. [2005] have introduced the *non-local means* approach to image denoising, where the prominence of self-similarities is used as a prior on natural images.⁹ Concretely, let us consider a noisy image written as a column vector \mathbf{y} in \mathbb{R}^n , and denote by $\mathbf{y}[i]$ the i -th pixel and, as usual, by \mathbf{y}_i the patch of size m centered on this pixel for some appropriate size m . This approach exploits the simple but very effective idea that two pixels associated with similar patches \mathbf{y}_i and \mathbf{y}_j should have similar values $\mathbf{y}[i]$ and $\mathbf{y}[j]$. Using \mathbf{y}_i as an explanatory variable for $\mathbf{y}[i]$ leads to the non-local means formulation, where the denoised pixel $\hat{\mathbf{x}}[i]$ is obtained by a weighted average:

$$\hat{\mathbf{x}}[i] = \sum_{j=1}^n \frac{K_h(\mathbf{y}_i - \mathbf{y}_j)}{\sum_{l=1}^n K_h(\mathbf{y}_i - \mathbf{y}_l)} \mathbf{y}[j], \quad (3.9)$$

and K_h is a Gaussian kernel of bandwidth h . The estimator (3.9) is in fact classical in the non-parametric statistics literature, even though it was never applied to image denoising with natural image patches before; it is usually called Nadaraya-Watson estimator [see Nadaraya, 1964, Watson, 1964, Wasserman, 2006].¹⁰

⁹This idea has in fact appeared in the literature in various guises and under different equivalent interpretations, *e.g.*, kernel density estimation, mean-shift iterations [Awate and Whitaker, 2006], diffusion processes on graphs [Szlám et al., 2007], and long-range random fields [Li and Huttenlocher, 2008].

¹⁰Interestingly, the same class of non-parametric estimators have also been used independently by Takeda et al. [2007], where the explanatory variables are not full

Later, several extensions have been proposed along this line of work. First attempts have focused on improving the computational speed [Mahmoudi and Sapiro, 2005], or on automatically adapting the kernel bandwidth h to the pixel of interest [Kervrann and Boulanger, 2006]. Extensions of the non-local principle for other tasks than denoising have also been developed, *e.g.*, for image upscaling [Glasner et al., 2009], or demosaicking [Buades et al., 2009]. Later, other estimators exploiting image self-similarities and based on nonparametric density estimation have been proposed and analyzed by Levin et al. [2012], Chatterjee and Milanfar [2012].

BM3D. Dabov et al. [2007a] proposed a patch-based procedure called BM3D for image denoising. The method exploits several ideas, including image self-similarities and sparse wavelet estimation, and provides outstanding results at a reasonable computational cost. Several years after it was developed, it remains considered as the state of the art. BM3D proceeds with the following pipeline for denoising an image \mathbf{y} :

- **block matching:** like non-local means, BM3D exploits self-similarity; it processes all patches by first matching them with similar ones in the noisy image \mathbf{y} , stacking them together into a 3D signal block;
- **3D filtering:** each block is denoised by using hard or soft-thresholding with a 3D orthogonal DCT dictionary, following the wavelet thresholding tradition;
- **patch averaging:** as in Section 3.1, a denoised image $\hat{\mathbf{x}}_0$ is obtained by averaging the estimates of the overlapping patches;¹¹
- **refinement with Wiener filtering:** the previous steps are refined by using the intermediate estimate $\hat{\mathbf{x}}_0$. First, block matching is applied to $\hat{\mathbf{x}}_0$ instead of \mathbf{y} in order to obtain a more reliable

image patches, but neighbor pixels. It results in Nadaraya-Watson estimators that exploit local information, as opposed to the non-local one of Buades et al. [2005].

¹¹More precisely, a weighted average is performed. We omit this detail here for simplicity, and refer to [Dabov et al., 2007a] for more details.

matching of noisy patches from \mathbf{y} ; weights are obtained for the 3D-filtering by using a Wiener filter. After a new patch averaging step, a final estimate $\hat{\mathbf{x}}_1$ is obtained.

A few additional heuristics are also used to further boost the denoising performance [see Dabov et al., 2007a]. Finally, the scheme has proven to be very efficient and gives substantially better results than regular non-local means. Later, it was improved by Dabov et al. [2009] by adding other components: (i) shape-adaptive patches that are non necessarily rectangular; (ii) 3D filtering with adaptive dictionaries obtained with PCA inside each block instead of using simple DCT.

Non-local sparse models. Since BM3D was successful in combining image self-similarities and wavelet/DCT denoising on blocks of similar patches, Mairal et al. [2009] have proposed to exploit further this idea and combine the non-local means principle with dictionary learning.

One motivation is that non-local means approach has proven to be effective in general, but it fails in some cases. In the extreme, when a patch does not look like any other one in the image, it is impossible to exploit self-similarities to estimate the corresponding pixel value. To some extent, sparse image models based on dictionary learning can handle such situations when the patch at hand admits a sparse decomposition, but they suffer from another drawback: similar patches sometimes admit very different estimates due to the potential instability of sparsity patterns, which can result in practice in noticeable reconstruction artifacts. The idea of the non-local sparse model is that *similar patches should admit similar sparse decompositions*. By enforcing this principle, one hopes to obtain more stable decompositions and subsequently a better estimation.

Concretely, let us define for each patch \mathbf{y}_i the set \mathcal{S}_i :

$$\mathcal{S}_i \triangleq \left\{ j = 1, \dots, n \text{ s.t. } \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \leq \xi \right\},$$

where ξ is some threshold. The block \mathcal{S}_i is essentially constructed by following the block matching step of BM3D. What differs is the way \mathcal{S}_i is subsequently processed. Whereas BM3D performs some wavelet filtering, the non-local sparse model jointly decomposes the patches onto

a previously learned dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$.

To encourage the decompositions to be similar, the patches from the set \mathcal{S}_i are forced to share a joint sparsity pattern—that is, they should share a common set of nonzero coefficients. Fortunately, this can be achieved with a *group-sparsity penalty*. Denoting by $\mathbf{A}_i \triangleq [\boldsymbol{\alpha}_{ij}]_{j \in \mathcal{S}_i}$ the matrix of coefficient in $\mathbb{R}^{p \times |\mathcal{S}_i|}$ corresponding to the group \mathcal{S}_i , the regularizer encourages the matrix \mathbf{A}_i to have a small number of non-zero rows, yielding a type of sparsity patterns that is illustrated in Figure 3.10.

Formally, the penalty can be written for a matrix \mathbf{A} in $\mathbb{R}^{p \times k}$ as

$$\Psi_q(\mathbf{A}) \triangleq \sum_{j=1}^p \|\boldsymbol{\alpha}^j\|_2^q,$$

where $\boldsymbol{\alpha}^j$ denotes the j -th row of \mathbf{A} . When $q = 1$, we obtain the group-Lasso norm already presented in Section 1.3. When $q = 0$, the penalty simply counts the number of non-zero rows in \mathbf{A} and thus plays the same role as the ℓ_0 -penalty in the context of group sparsity [Tropp et al., 2006]. Then, decomposing the patch \mathbf{y}_i with the penalty Ψ_q on the set \mathcal{S}_i amounts to solving

$$\min_{\mathbf{A}_i \in \mathbb{R}^{p \times |\mathcal{S}_i|}} \Psi_q(\mathbf{A}_i) \quad \text{s.t.} \quad \sum_{j \in \mathcal{S}_i} \|\mathbf{y}_j - \mathbf{D}\boldsymbol{\alpha}_{ij}\|_2^2 \leq \varepsilon_i.$$

As for the classical denoising approach based on dictionary learning presented in Section 3.1, the following ingredients are used to ultimately reconstruct the denoised image: (i) ε_i is chosen according to some effective heuristics; (ii) the non-convex penalty Ψ_0 with greedy algorithms [Tropp, 2004] is preferred to Ψ_1 for the final image reconstruction task once the dictionary has been learned; (iii) since the patches overlap, the final step of the algorithm averages the estimates of every pixel. A few additional heuristics are also used to improve the speed or quality of the algorithm [see Mairal et al., 2009], including a refinement step as in BM3D.

The non-local sparse principle was also successfully applied by Mairal et al. [2009] to image demosaicking (see Figure 3.4), and more generally to image interpolation [Romano et al., 2014].

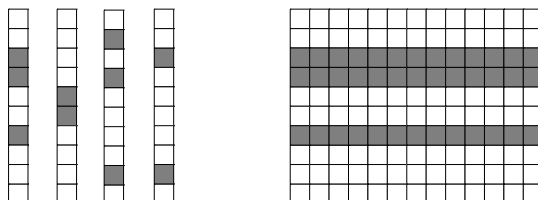


Figure 3.10: Sparsity vs. joint sparsity: Grey squares represents non-zeros values in vectors (left) or matrix (right).

Later, other variants of non-local sparse models have been proposed, notably the *centralized sparse representation* technique of Dong et al. [2011a, 2013]. There, input patches are clustered and an orthogonal dictionary is learned for each cluster with PCA. Then, patches are sparsely encoded on the PCA dictionaries, and the decompositions coefficients inside each cluster are encouraged to be close to a per-cluster mean representation. As shown in Sections 3.1 and 3.3, the performance achieved by non-local sparse models is in general competitive with the state of the art for different tasks, *e.g.*, denoising or demosaicking.

Gaussian mixture models. We have already mentioned in Section 2.7 that Gaussian mixture models have been successful for modeling natural image patches [Zoran and Weiss, 2011, Yu et al., 2012]. Not surprisingly, such approaches also lead to impressive results for image denoising and other image reconstructions tasks such as image upscaling [Yu et al., 2012]. Given a database of image patches, the model parameters are usually learned by using the EM-algorithm or one of its variants, and denoising can be achieved with classical estimators for probabilistic models, *e.g.*, maximum a posteriori estimation.

4

Sparse Coding for Visual Recognition

Because of its ability to model well natural image patches and automatically discover interpretable visual patterns, one may wonder whether dictionary learning can be useful for visual recognition or not. Patches are indeed often used as lowest-level features in image analysis pipelines. For instance, this is the case of popular approaches based on bags of words [Csurka et al., 2004, Lazebnik et al., 2006], which exploit image descriptors that are nothing else than pre-processed natural image patches [Lowe, 2004, Dalal and Triggs, 2005, Mikolajczyk and Schmid, 2005, Tola et al., 2010]. This is also the case of convolutional neural networks [LeCun et al., 1998a], whose first layer performs local convolutions on raw pixel values.

It is thus tempting to believe that modeling patches is an important step for automatic image understanding and that dictionary learning can be a key component of recognition architectures. The different approaches that we review in this section seem to confirm this fact, even though recognition requires other properties that are not originally provided by sparse coding models, such as invariance or stability of image representations to local perturbations.

In Section 4.1, we show how dictionary learning has been used for

image classification as an alternative to clustering techniques in bag-of-words models, before reviewing in Section 4.2 numerous variants that have been proposed in the literature. Then, we present different approaches where sparse models are used as a classification tool for face recognition (Section 4.3), patch classification, and edge detection (Section 4.4). Finally, we draw some links between dictionary learning and neural networks in Section 4.5 with backpropagation rules and networks for approximating sparse models, before reviewing a few additional computer vision applications in Section 4.6.

We show that dictionary learning has been successfully used in various guises for different recognition tasks.

4.1 A coding and pooling approach to image modeling

Following Boureau et al. [2010], we show in this section that interleaved steps of (nonlinear) feature coding and pooling are a very simple but common approach to image modeling in visual recognition tasks such as retrieval and categorization. It is, implicitly or explicitly, used in SIFT [Lowe, 2004], bags of features [Csurka et al., 2004, Sivic and Zisserman, 2003], HOG [Dalal and Triggs, 2005], the pyramid match kernel [Grauman and Darrell, 2005], spatial pyramids [Lazebnik et al., 2006], soft quantization [van Gemert et al., 2010], and most importantly for this presentation, sparse coding approaches to recognition [Yang et al., 2009, Boureau et al., 2010, Wang et al., 2010, Yang et al., 2010b].

We also show that the coding/pooling approach can often be intuitively justified in terms of (approximate) feature matching. We first consider *global pooling* approaches that discard all spatial information for every feature, before moving to *local pooling* that keeps part of it.

4.1.1 Global pooling

Bags of features and sparse coding. This paragraph largely follows Boureau et al. [2010]. Let us consider an image with its associated local image features, represented by (and identified with) a finite set $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ of n vectors in \mathbb{R}^m . For instance, each \mathbf{x}_i may

represent a non-rotationally invariant SIFT descriptor computed at location i in the image, but it may as well represent another feature type, *e.g.*, local binary pattern [Ojala et al., 2002] or a color histogram.

Let us also assume that we are given a *coding operator*, that is a function $\alpha : \mathbb{R}^m \rightarrow \mathcal{A}$ mapping features in \mathbb{R}^m onto the corresponding codes in some space \mathcal{A} , and a *pooling operator* β mapping finite subsets of \mathcal{A} onto elements of \mathbb{R}^p . The simplest version of the coding/pooling approach is to model an image represented by a set of descriptors \mathbf{X} in $\mathbb{R}^{m \times n}$ by the vector $\gamma(\mathbf{X}) = \beta[\alpha(\mathbf{X})]$, where $\alpha(\mathbf{X})$ denotes the set of points $[\alpha(\mathbf{x}_i)]_{i=1}^n$ in \mathcal{A}^n . We dub this approach *global pooling* because it summarizes the codes of *all* features in the image by a single vector in \mathbb{R}^p .

The best known example of global pooling is the bag-of-features model inherited from text processing and popularized in computer vision by Sivic and Zisserman [2003] and Csurka et al. [2004]: let us consider a quantizer $\xi : \mathbb{R}^m \rightarrow \{1, \dots, p\}$, constructed (for example) using the algorithm K-means on the elements of \mathbf{X} . We can use ξ to encode each feature \mathbf{x} as the binary vector $\alpha(\mathbf{x})$ of dimension p with all zero entries except for entry number $\xi(\mathbf{x})$, which is equal to one. Note that the code space \mathcal{A} is $\{0, 1\}^p$ in this case.

We can now use *mean pooling* (also called *average pooling*) to summarize the feature codes $\alpha(\mathbf{x}_i)$ of the whole image by their average—that is, a single vector of dimension p :¹

$$\gamma(\mathbf{X}) = \beta[\alpha(\mathbf{X})], \quad \text{where} \quad \beta[\alpha(\mathbf{X})] = \frac{1}{n} \sum_{i=1}^n \alpha(\mathbf{x}_i).$$

The vector $\gamma(\mathbf{X})$ is called the bag of features associated with the image.

Two images respectively represented by the local descriptors $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$ and $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_{n'}]$ in $\mathbb{R}^{m \times n'}$ can now be compared by computing the dot product of their bags of features, that is, their similarity is computed as $s(\mathbf{X}, \mathbf{Y}) = \gamma(\mathbf{X})^\top \gamma(\mathbf{Y})$. This approach

¹Note that the tf-idf renormalization common in image retrieval [Sivic and Zisserman, 2003] can be achieved by reweighting the components of $\gamma(\mathbf{X})$: let γ_i denote the i -th component of $\gamma(\mathbf{X})$, *i.e.*, the frequency of the corresponding bin in \mathbf{X} , replace γ_i by $\gamma_i \log(N/n\gamma_i)$, where N is the total number of features in the database being queried.

is commonly used in both image retrieval (where s is used for ranking) [Sivic and Zisserman, 2003] and categorization (where s is used as a kernel) [Csurka et al., 2004]. In both cases, the similarity function can be computed efficiently due to the fact that bags of features are typically very sparse for large values of p .

When the quantizer ξ is obtained using the K-means algorithm, the centroids \mathbf{d}_j ($j = 1, \dots, p$) of the corresponding clusters can be seen as the elements (“atoms”, or “words”) of some visual dictionary, and the codes $\boldsymbol{\alpha}(\mathbf{x})$ tell us what atom a given feature is associated with.

This approach readily generalizes to other types of feature coding. In particular, it is possible to replace the vector-quantized binary codes of the features \mathbf{x} in \mathbf{X} by their sparse code relative to some (given or learned) dictionary \mathbf{D} [Yang et al., 2009, Boureau et al., 2010], that is

$$\boldsymbol{\alpha}(\mathbf{x}) \in \arg \min_{\boldsymbol{\alpha}' \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}'\|_2^2 + \lambda \|\boldsymbol{\alpha}'\|_1,$$

then use mean pooling as before.

Feature matching interpretation. Here, we follow Jégou et al. [2010], and show that, in the case of bags of features constructed with a quantifier, the similarity function $s(\mathbf{X}, \mathbf{Y})$ can be justified intuitively by the fact that it measures the proportion of pairs of features in \mathbf{X} and \mathbf{Y} that match each other [Jégou et al., 2010]. Maximizing $s(\mathbf{X}, \mathbf{Y})$ can thus be thought of as some (approximate) feature matching process. Indeed, given a quantizer ξ , let us define the functions $\delta_j : \mathbb{R}^m \rightarrow \{0, 1\}$ ($j = 1, \dots, p$) and $\delta : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \{0, 1\}$ by

$$\delta_j(\mathbf{x}) = \begin{cases} 1 & \text{when } \xi(\mathbf{x}) = j, \\ 0 & \text{otherwise,} \end{cases} \quad \text{and} \quad \delta(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{when } \xi(\mathbf{y}) = \xi(\mathbf{x}), \\ 0 & \text{otherwise.} \end{cases} \quad (4.1)$$

We can now rewrite our similarity function as

$$s(\mathbf{X}, \mathbf{Y}) = \sum_{j=1}^p \left[\frac{1}{n} \sum_{i=1}^n \delta_j(\mathbf{x}_i) \right] \left[\frac{1}{n'} \sum_{i'=1}^{n'} \delta_j(\mathbf{y}_{i'}) \right] = \frac{1}{n n'} \sum_{i=1}^n \sum_{i'=1}^{n'} \delta(\mathbf{x}_i, \mathbf{y}_{i'}),$$

which indeed measures the fraction of pairs of features in the two images that are deemed to match when they admit the same code.

When the quantizer ξ is obtained by running the K-means algorithm, $\delta(\mathbf{x}, \mathbf{y})$ is nonzero when the features \mathbf{x} and \mathbf{y} both lie in the same cell of the Voronoi diagram of the p centroids of the corresponding clusters. Related approaches explicitly phrased in term of feature correspondences have been used in both image matching and retrieval using some nearest-neighbor (or NN) decision rule such as ε -search or K -NN. Concretely, one measures again the fraction of matching features as

$$s(\mathbf{X}, \mathbf{Y}) = \frac{1}{n n'} \sum_{i=1}^n \sum_{i'=1}^{n'} \delta(\mathbf{x}_i, \mathbf{y}_{i'}),$$

but, this time the function δ is defined by

$$\delta(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{when } \|\mathbf{y} - \mathbf{x}\|_2 < \varepsilon, \\ 0 & \text{otherwise,} \end{cases} \quad (4.2)$$

in the ε -search case, or, in the K -NN case, by

$$\delta(\mathbf{x}, \mathbf{y}) = \begin{cases} 1 & \text{when } \mathbf{y} \text{ is one of the } K \text{ nearest neighbors of } \mathbf{x}, \\ 1 & \text{when } \mathbf{x} \text{ is one of the } K \text{ nearest neighbors of } \mathbf{y}, \\ 0 & \text{otherwise.} \end{cases} \quad (4.3)$$

This approach is the basis of the classical image retrieval and image matching techniques of Schmid and Mohr [1997] and Lowe [2004]. As in the case of bags of features [Sivic and Zisserman, 2003], the matching step is often complemented by some geometrical verification stage in practice. Similar schemes have also been used in image categorization [Wallraven et al., 2003].

It should be noted that the coding/pooling models also admits a feature-matching interpretation when images are encoded using sparse coding instead of a quantifier: we can write in this case the similarity $s(\mathbf{X}, \mathbf{Y})$ of two images, respectively represented by the sets of descriptors \mathbf{X} and \mathbf{Y} , as

$$s(\mathbf{X}, \mathbf{Y}) = \sum_{j=1}^p \left[\frac{1}{n} \sum_{i=1}^n \alpha_j(\mathbf{x}_i) \right] \left[\frac{1}{n'} \sum_{i'=1}^{n'} \alpha_j(\mathbf{y}_{i'}) \right] = \frac{1}{n n'} \sum_{i=1}^n \sum_{i'=1}^{n'} \delta(\mathbf{x}_i, \mathbf{y}_{i'}),$$

where $\boldsymbol{\alpha}(\mathbf{x}) = [\alpha_1(\mathbf{x}), \dots, \alpha_p(\mathbf{x})]$ and $\delta(\mathbf{x}, \mathbf{y}) = \boldsymbol{\alpha}(\mathbf{x})^\top \boldsymbol{\alpha}(\mathbf{y})$. Note that by analogy, $\delta(\mathbf{x}, \mathbf{y})$ can be rewritten in a similar form to the previous

cases (4.1), (4.3), or (4.2):

$$\delta(\mathbf{x}, \mathbf{y}) = \begin{cases} \boldsymbol{\alpha}(\mathbf{x})^\top \boldsymbol{\alpha}(\mathbf{y}) & \text{when } \exists j, \alpha_j(\mathbf{x})\alpha_j(\mathbf{y}) \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$

The function δ encodes in this case the fact that two feature vectors are deemed to match when they both have at least one nonzero code element for some dictionary atom, the contribution of each matching feature pair being proportional to the dot product of their codes.

Max pooling. The construction of the vector $\boldsymbol{\gamma}(\mathbf{X})$ associated with an image can be seen as pooling the code vectors $[\boldsymbol{\alpha}(\mathbf{x}_i)]_{i=1}^n$ associated with the image features and summarizing them by a single vector $\boldsymbol{\gamma}(\mathbf{X})$ which is their average. As noted before, this process is thus often dubbed *mean*, or *average pooling*. A simple variant is *sum pooling*, where the vectors $\boldsymbol{\gamma}(\mathbf{X})$ are not normalized and simply sum up the vectors $\boldsymbol{\alpha}(\mathbf{x}_i)$.

It has proven useful in several classification tasks to replace mean pooling by *max pooling* [Serre et al., 2005], where

$$\gamma_j(\mathbf{X}) = \max_{i=1, \dots, n} \alpha_j(\mathbf{x}_i),$$

and $\boldsymbol{\gamma}(\mathbf{X}) = [\gamma_1(\mathbf{X}), \dots, \gamma_p(\mathbf{X})]$. In the case of bags of features, the vectors $\boldsymbol{\alpha}(\mathbf{x}_i)$ are binary, and computing the similarity $s(\mathbf{X}, \mathbf{Y})$ amounts to counting the number of bins in the quantization associated with at least one feature in each image. Intuitively, this may be justified in an application such as Video Google Sivic and Zisserman [2003], where the query image (typically a box drawn around an object of interest such as a tie or a plate) is normally much smaller than the database images (typically depicting an entire scene) so that features occurring often in the scene are not given too much importance.

Max pooling is also used with sparse coding. This is a priori problematic because of the sign ambiguity of dictionary learning discussed in Section 1.6. Indeed, as noted by Murray and Perronnin [2014], max pooling should only be used (or perhaps more accurately, is only intuitively justified) when the individual vector entries encode a strength of association between a descriptor and a codeword, and are thus all

positive. This is the case for bags of features but not for sparse coding, with max pooling resulting potentially in 2^p different pooled features. The construction of Section 1.6 can be used to remove this ambiguity, but does not provide an intuitive justification of max pooling in this setting. It may then make sense to use a variant of dictionary learning that enforces non-negativity constraints on the sparse codes, or to duplicate the entries of $\alpha(\mathbf{x})$ into negative and positive components, resulting in a vector of size $2p$, before applying max pooling. Specifically, each entry $\alpha_j(\mathbf{x})$ will be duplicated into two values $\max(\alpha_j(\mathbf{x}), 0)$ and $\max(-\alpha_j(\mathbf{x}), 0)$.

Finally, other nonlinear pooling techniques have also been studied in the literature [see Koniusz et al., 2013]; for simplicity, we restrict our presentation to the average and max-pooling strategies, which are the most popular ones.

4.1.2 Local pooling

Spatial pooling: spatial pyramids and their cousins The global coding/pooling approach yields “orderless” image models Koenderink and Van Doorn [1999], where all spatial information has been discarded. While this affords a great deal of robustness (including a total invariance to image transformations *as long as* they leave the local features invariant),² it seems wasteful to discard *all* spatial information.

The spatial pyramid of Lazebnik et al. [2006] addresses this problem by overlaying a coarse pyramid with L levels over the image (the HOG model of Dalal and Triggs [2005] is based on a similar idea using a coarse grid instead of a pyramid). There are 2^{2l} cells at level l ($l = 0, \dots, L - 1$) of the pyramid, and the features falling in each cell of the pyramid are binned separately. Let us define by \mathcal{S}_{kl} the indices of features in $\{1, \dots, n\}$ falling into cell number k at level l ; the spatial pyramid image descriptor $\gamma(\mathbf{X})$ is obtained by concatenating the

²One should keep in mind that being invariant to within-image transformations does not imply being invariant to, say, rotations in depth since (at least) some features will become occluded, or will be revealed, as the observed object rotates relative to the camera.

unnormalized histograms (sum pooling) $\gamma_{kl}(\mathbf{X}) = \sum_{i \in \mathcal{S}_{kl}} \alpha(\mathbf{x}_i)$ associated with all cells at all levels of the pyramid to form a vector in \mathbb{R}^{pd} , where $d = 1 + \dots + d_{L-1} = (2^{2L} - 1)/3$, and $d_l = 2^{2l}$ for $l = 0, \dots, L-1$.

Note that the case $L = 1$ corresponds to the global pooling model, as used in its bag of features or sparse coding instances, except for the fact that it uses *sum pooling* instead of mean or max pooling, that is, the sum of the features is used instead of their mean or max value. As explained in the next section, the spatial pyramids associated with two images can be compared using their dot product or the *pyramid matching kernel* of Grauman and Darrell [2005], and they can be interpreted in terms of feature matching in both cases.

Feature matching interpretation. Given again two set of descriptors \mathbf{X} and \mathbf{Y} and a spatial pyramid structure leading to respective sets of indices \mathcal{S}_{kl} and \mathcal{S}'_{kl} for the cell k at level l , we can compute the similarity

$$s(X, Y) = \gamma(\mathbf{X})^\top \gamma(\mathbf{Y}) = \sum_{l=0}^{L-1} \sum_{k=1}^{d_l} \sum_{i \in \mathcal{S}_{kl}} \sum_{i' \in \mathcal{S}'_{kl}} \alpha(\mathbf{x}_i)^\top \alpha(\mathbf{y}_{i'}).$$

This similarity function measures the number of matches between the two images, where matches are found at different spatial scales as pairs of features falling in the same cell.

As shown by Lazebnik et al. [2006], an alternative is provided by the *pyramid match kernel*, defined as follows. Let us first define the *histogram intersection* function as

$$I_l(\mathbf{X}, \mathbf{Y}) = \sum_{k=1}^{d_l} \min[\gamma_{kl}(\mathbf{X}), \gamma_{kl}(\mathbf{Y})],$$

where $\gamma_{kl}(\mathbf{Y}) = \sum_{i \in \mathcal{S}'_{kl}} \alpha(\mathbf{y}_i)$, $\gamma_{kl}(\mathbf{X}) = \sum_{i \in \mathcal{S}_{kl}} \alpha(\mathbf{x}_i)$, and where the min operator is applied componentwise. The j -th coordinate of $I_l(\mathbf{X}, \mathbf{Y})$ measures the number of matches between \mathbf{X} and \mathbf{Y} corresponding to the atom \mathbf{d}_j of the dictionary, measured as the number of points from \mathbf{X} and \mathbf{Y} that have nonzero codes, and fall in the same cell. These points match in the code space because they have a nonzero element in the

same spot, and they also match spatially, because they fall into the same cell.

By construction, the matches found at level l of the pyramid also include the matches found at level $l+1$. The number of new matches found at level l is thus $I_l(\mathbf{X}, \mathbf{Y}) - I_{l+1}(\mathbf{X}, \mathbf{Y})$ for $l = 0, \dots, L-1$. This suggests the following *pyramid match kernel* [Grauman and Darrell, 2005] to compare two images:

$$\begin{aligned} K(\mathbf{X}, \mathbf{Y}) &= I_L(\mathbf{X}, \mathbf{Y}) + \sum_{l=0}^{L-1} \frac{1}{2^{L-l}} [I_l(\mathbf{X}, \mathbf{Y}) - I_{l+1}(\mathbf{X}, \mathbf{Y})] \\ &= \frac{1}{2^L} I_0(\mathbf{X}, \mathbf{Y}) + \sum_{l=0}^{L-1} \frac{1}{2^{L-l+1}} I_l(\mathbf{X}, \mathbf{Y}), \end{aligned}$$

where the weight $\frac{1}{2^{L-l}}$, which is inversely proportional to the cell width at level l is used to penalize matches found in larger cells because they involve increasingly dissimilar features. It is easily shown that K is a positive-definite kernel, which allows using the machinery of kernel methods and reproducing kernel Hilbert spaces [Shawe-Taylor and Cristianini, 2004].

Like the histogram intersection function, the pyramid match kernel can be interpreted as measuring the number of matches between \mathbf{X} and \mathbf{Y} in both feature space and the spatial domain. Contrary to the previous cases discussed in this presentation, where all pairs of matching features were counted, a point of \mathbf{X} may only match a unique point of \mathbf{Y} in this case.

Feature space pooling. The pyramid match kernel was originally proposed by Grauman and Darrell [2005] as a method for matching images (or equivalently counting approximate correspondences) in *feature space*, without retaining any spatial information. With spatial pyramids, Lazebnik et al. [2006] argued that it might be better suited to encoding spatial information in the two-dimensional image, using instead traditional vector quantization based on the K-means algorithm to handle matching in the high-dimensional feature space.

The two approaches can be thought of as performing a *local* form of pooling, over bins defined in feature space in Grauman and Darrell

[2005], or over cells defined in image space in Lazebnik et al. [2006]. In practice, the latter method usually gives better results for image categorization, but it may be interesting to combine the two when the feature codes are obtained using sparse coding. This is precisely what Boureau et al. [2011] have proposed to do [see also Gao et al., 2010, 2013, Wang et al., 2010, Yang et al., 2012b, Zhou et al., 2010, Koniusz and Mikolajczyk, 2011, for related work that will be partly detailed in the next section]. In this approach, a dictionary is learned on the training data, and the corresponding sparse codes are then clustered using K-means. Let \mathcal{S}_{klj} now denote the index set of features that fall in cell number k of the image at level l of the pyramid, and in the Voronoi cell number j of the sparse code space, the unnormalized histogram $\gamma_{kl}(\mathbf{X})$ of spatial pyramid is simply replaced by $\gamma_{klj}(\mathbf{X}) = \sum_{i \in \mathcal{S}_{klj}} \alpha(\mathbf{x}_i)$, resulting in a final descriptor $\gamma(\mathbf{X})$ of size pdK where K is the number of centroids used by the K-means algorithm. Combined with max pooling, this method indeed gives very good results for image categorization on standard benchmarks such as Caltech 101.

4.2 The botany of sparse feature coding

To address the limitations of the original dictionary learning formulation of Olshausen and Field [1996, 1997] for feature encoding, many different variants have been proposed in the literature. In this section, we go through a few of them that have gained a significant amount of attention and have been reported to improve upon the original codebook learning approach presented in the previous section.

Local coordinate coding. Yu et al. [2009] have proposed a significantly different view of the dictionary learning problem than the one presented so far in this monograph, where the learned dictionary elements are interpreted as anchor points on a nonlinear smooth manifold.³ More precisely, these anchor points define a local coordinate system to represent data points, *e.g.*, natural image patches or local

³Note that the definition of “manifold” from Yu et al. [2009] slightly differs from the traditional one. More precisely, they define a “manifold” \mathcal{M} as a subset of \mathbb{R}^m such there exists a constant C and “dimension” p such that for all \mathbf{x} in \mathcal{M} , there

descriptors. To make this interpretation relevant, Yu et al. [2009] have introduced the concept of *locality* for the dictionary elements, encouraging them to be close to data points. In practice, the resulting formulation, called “local coordinate coding”, consists of optimizing the following cost function:

$$\min_{\mathbf{D} \in \mathbb{R}^{m \times p}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \sum_{j=1}^p \|\mathbf{x}_i - \mathbf{d}_j\|_2^2 |\boldsymbol{\alpha}_i[j]| \right), \quad (4.4)$$

where, as usual, $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$ is the database of training signals, $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n]$ is the matrix of sparse coefficients, and $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ is the dictionary. In contrast to the classical dictionary learning formulation, the regularization function is a weighted ℓ_1 -norm, where the quadratic weights $\|\mathbf{x}_i - \mathbf{d}_j\|_2^2$ encourage the selection of dictionary elements that are close to the signal \mathbf{x}_i in Euclidean norm. Optimizing (4.4), meaning finding a stationary point since the problem is nonconvex, can be achieved with alternate minimization between \mathbf{D} and \mathbf{A} , which is a classical strategy for dictionary learning (see Section 5).

Note that the link between (4.4) and the manifold assumption is not obvious at first sight. In fact, Yu et al. [2009] show that (4.4) can be interpreted as a practical heuristic for a slightly different formulation that has precise theoretical guarantees, namely an approximation bound for modeling the supposedly existing manifold. Later, the locality principle was revisited by Yu and Zhang [2010] with an improved approximation scheme, and by Wang et al. [2010] with a practical feature encoding scheme.

The latter approach, dubbed “locality-constrained linear coding (LLC)”, exploits the locality principle in a simple way to select a predefined number $K < p$ of non-zero coefficients for the codes $\boldsymbol{\alpha}_i$ instead of using the ℓ_1 -regularization. Given a fixed dictionary \mathbf{D} and a signal \mathbf{x}_i , the LLC scheme forces the support of $\boldsymbol{\alpha}_i$ to correspond

exists p vectors $\mathbf{D}(\mathbf{x}) = [\mathbf{d}_1(\mathbf{x}), \dots, \mathbf{d}_p(\mathbf{x})]$ in $\mathbb{R}^{m \times p}$ such that for all \mathbf{x}' in \mathbb{R}^m ,

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{x}' - \mathbf{x} - \mathbf{D}(\mathbf{x})\boldsymbol{\alpha}\|_2 \leq C \|\mathbf{x}' - \mathbf{x}\|_2^2.$$

to the K -nearest dictionary elements to \mathbf{x}_i . The value of the non-zero coefficients are then computed in analytical form without having to solve a sparse regression problem. Recent reviews and benchmarks have shown that such a coding scheme is competitive for image classification tasks [Sánchez et al., 2013].

Laplacian sparse coding. We have seen in Section 3.8 with the non-local sparse models that a successful idea for image restoration is to encourage similar signals to share similar sparse decomposition patterns. For visual recognition, Gao et al. [2010, 2013] have shown that such a principle can be also useful. Given a training set $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ of local descriptors, their formulation called “Laplacian sparse coding” consists of minimizing a regularized dictionary learning objective function

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \left(\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1 \right) + \mu \sum_{i,j} W_{ij} \|\boldsymbol{\alpha}_i - \boldsymbol{\alpha}_j\|_2^2,$$

where $\mathbf{W} = [W_{ij}]$ in $\mathbb{R}^{n \times n}$ is a similarity matrix; W_{ij} should be typically large when \mathbf{x}_i and \mathbf{x}_j are similar according to some appropriate metric, encouraging therefore $\boldsymbol{\alpha}_i$ to be close to $\boldsymbol{\alpha}_j$ in Euclidean norm. The terminology of “Laplacian” comes from the machine learning literature about semi-supervised learning where model variables sit on a graph [Belkin and Niyogi, 2003, 2004] and where such regularization functions are used. Gao et al. [2010, 2013] further discuss the choice of a good similarity matrix \mathbf{W} when the signals \mathbf{x}_i are local descriptors, and experimentally show that the additional regularization yields better classification results than the traditional sparse coding approach of Yang et al. [2009].

Convolutional sparse coding. In Section 2.6, we have presented the convolutional sparse coding model applied to natural images without any concrete application. Zeiler et al. [2011] have applied such a formulation for visual recognition with a multilayer scheme inspired from convolutional neural networks and from other hierarchical models [Serre et al., 2005].

For each layer of the network, the formulation leverages the concept of “feature maps”—that is, a set of two-dimensional spatial maps that carry coefficients corresponding to a particular feature type at different locations. A feature map \mathbf{x} has typically three dimensions, *e.g.*, \mathbf{x} is in $\mathbb{R}^{\sqrt{l} \times \sqrt{l} \times p}$ for a square feature map of spatial size $\sqrt{l} \times \sqrt{l}$ and p different feature types, but it is sometimes more convenient to represent \mathbf{x} as a vector of size pl . It is also convenient to consider three-dimensional spatial patches in the feature maps, *e.g.*, a patch of size $\sqrt{e} \times \sqrt{e} \times p$ that contains all information across feature maps from a spatial window of size $\sqrt{e} \times \sqrt{e}$ centered at a particular location.

Assuming that the parameters of the network have been already learned, one layer of the hierarchical model processes input feature maps produced by the previous layer (the first input feature map at the bottom of the hierarchy being the image itself), and produces an output feature map. Each layer performs successively two operations, which are illustrated in Figure 4.1:

1. **convolutional sparse coding**: the input feature map \mathbf{x} is encoded by using the convolutional sparse coding formulation already presented in Section 2.6, which we recall now. Assuming that a dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$ has been previously learned, the feature map \mathbf{x} —represented as a vector in \mathbb{R}^l here—is encoded as follows:

$$\min_{\mathbf{A} \in \mathbb{R}^{p \times l}} \frac{1}{2} \left\| \mathbf{x} - \sum_{k=1}^l \mathbf{R}_k^\top \mathbf{D} \boldsymbol{\alpha}_k \right\|_2^2 + \lambda \sum_{k=1}^l \|\boldsymbol{\alpha}_k\|_1,$$

where \mathbf{R}_k is the linear operator that extracts the patch centered at the k -th location from \mathbf{x} and \mathbf{R}_k^\top positions a small patch at a location k in a larger feature map, using a similar notation as in (3.3). This operation produces sparse codes $\mathbf{A} = [\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_l]$ in $\mathbb{R}^{p \times l}$, with p different coefficients for each position k . Then, the matrix \mathbf{A} can be interpreted as an intermediate feature map of size $\sqrt{l} \times \sqrt{l}$ with p channels, as shown in Figure 4.1.

2. **three-dimensional max pooling**: The spatial resolution and the number of channels of \mathbf{A} is typically reduced by 2 with a max-

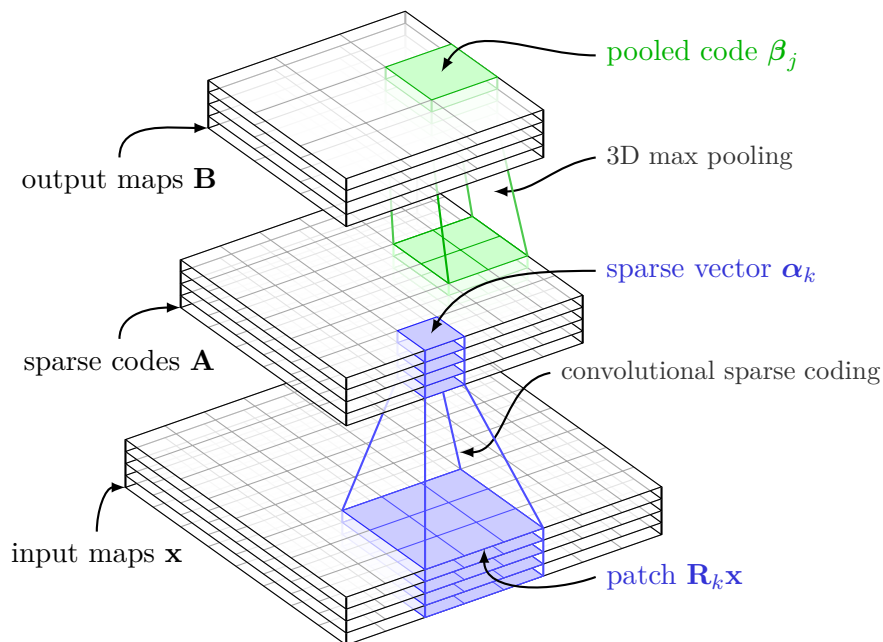


Figure 4.1: Illustration representing one layer of the hierarchical convolutional sparse coding model of Zeiler et al. [2011]. The vector \mathbf{x} denotes the input feature maps, which are encoded into sparse codes $\mathbf{A} = [\alpha_k]_{k=1}^l$ with the convolutional sparse coding formulation. The output feature maps $\mathbf{B} = [\beta_j]_{j=1}^{l'}$ are obtained after a three-dimensional max pooling step, reducing the spatial resolution and the number of maps from p to p' .

pooling operation already presented in Section 4.1 with three-dimensional pooling regions.

Separable sparse coding. A variant of the convolutional sparse coding model has also been investigated by Rigamonti et al. [2013], where the dictionary elements are spatially separable, or partially separable. More precisely, since a dictionary element \mathbf{d}_j in \mathbb{R}^m represents a two-dimensional square patch, it is possible to reorganize its entries as a matrix $\text{Mat}(\mathbf{d}_j)$ in $\mathbb{R}^{\sqrt{m} \times \sqrt{m}}$. Then, a dictionary element \mathbf{d}_j is said to be separable when there exist some vectors \mathbf{u}_j and \mathbf{v}_j in $\mathbb{R}^{\sqrt{m}}$ such that $\text{Mat}(\mathbf{d}_j) = \mathbf{u}_j \mathbf{v}_j^\top$. In other words, separability simply means that the matrix $\text{Mat}(\mathbf{d}_j)$ is rank one.

The motivation for looking for separable dictionary elements is to accelerate the computation of the inner products $\mathbf{d}_j^\top \mathbf{R}_k \mathbf{x}$ for all overlapping patches $\mathbf{R}_k \mathbf{x}$ from an input image \mathbf{x} in \mathbb{R}^l . This step indeed dominates the computational cost of reconstruction algorithms for the convolutional sparse coding model, and thus it is critical to make it efficient. The key observation is that computing all products $\mathbf{d}_j^\top \mathbf{R}_k \mathbf{x}$ is equivalent to performing a convolution of the two-dimensional filter $\text{Mat}(\mathbf{d}_j)$ on the input image \mathbf{x} , which naively requires $O(ml)$ operations.⁴ In the separable case, we remark that we have $\mathbf{d}_j^\top \mathbf{R}_k \mathbf{x} = \text{trace}(\text{Mat}(\mathbf{d}_j)^\top \text{Mat}(\mathbf{R}_k \mathbf{x})) = \text{trace}(\mathbf{v}_j \mathbf{u}_j^\top \text{Mat}(\mathbf{R}_k \mathbf{x})) = \mathbf{u}_j^\top \text{Mat}(\mathbf{R}_k \mathbf{x}) \mathbf{v}_j$. It is then easy to see that the inner products can be obtained by convolving twice the input image \mathbf{x} : first with the filter \mathbf{u}_j along the vertical direction, then with \mathbf{v}_j^\top along the horizontal direction. As a result, the complexity drops to $O(\sqrt{ml})$ operations.

Rigamonti et al. [2013] further study these filters for three-dimensional data, where the gain in complexity is even more important than in the two-dimensional case. They also show how to learn partially separable filters by using the following regularization function for the dictionary:

$$\varphi(\mathbf{D}) \triangleq \sum_{j=1}^p \|\text{Mat}(\mathbf{d}_j)\|_*,$$

where $\|\cdot\|_*$ is the trace norm presented in Section 1.3. As a result, the penalty φ encourages the matrices $\text{Mat}(\mathbf{d}_j)$ to be low-rank instead of exactly rank one, meaning that they can be expressed by a sum of a few rank-one matrices.

Another variant of “separable dictionary learning” has also been proposed at the same time by Hawe et al. [2013] in a non-convolutional setting, where the dictionary is assumed to be the Kronecker product of two smaller dictionaries:

$$\mathbf{D} = \mathbf{B} \otimes \mathbf{C},$$

⁴Note that the fast Fourier transform (FFT) could be used to reduce this theoretical complexity. Unfortunately, the FFT induces a computational overhead that make it inefficient in practice when convolving a large image with a small filter, which is the case here.

where \mathbf{B} is in $\mathbb{R}^{m_b \times p_b}$ and \mathbf{C} is in $\mathbb{R}^{m_c \times p_c}$. In other words, the matrix \mathbf{D} is in $\mathbb{R}^{m_b m_c \times p_b p_c}$ and is a block matrix with $m_b \times p_b$ blocks, such that the i, j -th block is $\mathbf{B}[i, j]\mathbf{C}$. Similar to Rigamonti et al. [2013], the motivation of Hawe et al. [2013] is to make the matrix-vector multiplications $\mathbf{D}\boldsymbol{\alpha}$ or $\mathbf{D}^\top \mathbf{x}$ more efficient by exploiting the properties of Kronecker products [see Golub and Van Loan, 2012]. However, it is worth noticing that the concepts of “separability” used by Rigamonti et al. [2013] and Hawe et al. [2013] are not equivalent to each other.

Multipath sparse coding and hierarchical matching pursuit. Finally, another model called hierarchical matching pursuit and introduced by Bo et al. [2011, 2013] has recently obtained state-of-the-art results in various recognition tasks. Similar to the convolutional approach of Zeiler et al. [2011], the scheme of Bo et al. [2011] produces for each image a sequence of feature maps organized in a multilayer fashion, but the coding scheme is somewhat simpler than in the convolutional sparse coding model. In a nutshell, the main differences are the following:

- **sparse coding of independent patches:** instead of using the convolutional model, all patches from an input feature map are coded independently;
- **orthogonal matching pursuit:** instead of using the ℓ_1 -regularization, the coding algorithm is a greedy method, which is presented in Section 5.1;
- **contrast normalization:** after each pooling step, patches are contrast-normalized, following the procedure described in Section 2.1.

Later, the multilayer scheme was improved by Bo et al. [2013] by combining several networks that have a different number of layers and thus that have different invariant properties. A simpler version of this scheme with a single layer has also been shown to be effective for replacing low-level features such as histograms of gradients [Dalal and Triggs, 2005] for object detection tasks [Ren and Ramanan, 2013].

4.3 Face recognition

One important success of sparse estimation for classification tasks is face recognition. Wright et al. [2009] have indeed proposed a state-of-the-art classifier for this task; as we will see, however, even though it was first evaluated on face datasets, the classifier is generic and can be applied to other modalities. More specifically, let us first consider a set of training signals from k different classes, represented by the matrix $\mathbf{D} = [\mathbf{D}_1, \dots, \mathbf{D}_k]$ in $\mathbb{R}^{m \times p}$ where the columns of each sub-matrix \mathbf{D}_j are signals from the j -th class. In order to classify a new test signal \mathbf{x} in \mathbb{R}^m , Wright et al. [2009] have proposed to sparsely decompose \mathbf{x} onto the matrix \mathbf{D} , and then measure which one of the sub-matrix \mathbf{D}_j is the most “used” in the decomposition. The precise procedure is given in Algorithm 1.

Interestingly, this approach called “sparse-representation based classifier (SRC)” is related to other non-parametric machine learning techniques such as nearest neighbor classifiers that look for the most similar training samples to \mathbf{x} before choosing the class label with a voting scheme [see, *e.g.*, Hastie et al., 2009], or methods looking for the nearest subspace [Naseem et al., 2010], which project the data \mathbf{x} onto each span of the submatrices \mathbf{D}_j , before selecting the subspace that best reconstructs the input data.

For the problem of face recognition, each matrix \mathbf{D}_j contains face samples of the same subject, and the method of Wright et al. [2009] has gained a lot of popularity. However, this approach may suffer from several issues, which have been addressed by various extensions:

- **occlusion:** face images in realistic environments often contain unwanted sources of variations. Subjects may indeed wear clothes, scarves, glasses, which occlude part of the face, and which are not necessarily present in the training data. To improve the robustness of the formulation, Wright et al. [2009] have introduced an auxiliary variable \mathbf{e} in \mathbb{R}^p whose purpose is to model occluded parts. Then, Eq. (4.5) is replaced by the following linear program:

$$(\boldsymbol{\alpha}^*, \mathbf{e}^*) \in \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p, \mathbf{e} \in \mathbb{R}^m} [\|\boldsymbol{\alpha}\|_1 + \|\mathbf{e}\|_1 \text{ s.t. } \mathbf{x} = \mathbf{D}\boldsymbol{\alpha} + \mathbf{e}],$$

Algorithm 1 Sparse-representation based Classifier (SRC).

Require: training instances $\mathbf{D} = [\mathbf{D}_1, \dots, \mathbf{D}_k]$ in $\mathbb{R}^{m \times p}$, each \mathbf{D}_j represents signals from class j ; One test signal \mathbf{x} in \mathbb{R}^m .

- 1: **normalization:** make each column of \mathbf{X} of unit ℓ_2 -norm;
- 2: **sparse decomposition:**

$$\boldsymbol{\alpha}^* \in \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} [\|\boldsymbol{\alpha}\|_1 \text{ s.t. } \mathbf{x} = \mathbf{D}\boldsymbol{\alpha}]; \quad (4.5)$$

- 3: **classification:**

$$\hat{j} \in \arg \min_{j \in \{1, \dots, k\}} \|\mathbf{x} - \mathbf{D}_j \delta_j(\boldsymbol{\alpha}^*)\|_2^2, \quad (4.6)$$

where δ_j selects the entries of $\boldsymbol{\alpha}^*$ corresponding to the class j ; in other words, $\mathbf{D}\boldsymbol{\alpha}^* = \sum_{j=1}^k \mathbf{D}_j \delta_j(\boldsymbol{\alpha}^*)$;

- 4: **return** the estimated class \hat{j} for the signal \mathbf{x} .
-

where the ℓ_1 -norm encourages the variable \mathbf{e} to be sparse. Thus, one hopes the variable \mathbf{e} to have non-zero coefficients corresponding to occluded areas that are difficult to reconstruct with the training data \mathbf{D} . The selection rule (4.6) needs also to be modified, and becomes

$$\hat{j} \in \arg \min_{j \in \{1, \dots, k\}} \|\mathbf{x} - \mathbf{D}_j \delta_j(\boldsymbol{\alpha}^*) - \mathbf{e}^*\|_2^2.$$

- **image misalignment:** the main assumption of the face recognition system of Wright et al. [2009] is that a new test image can be sparsely decomposed as a linear combination of training images. Of course, such an assumption does not hold when the training images are not perfectly aligned—that is, when the face features for a subject are not localized at the same positions across training samples. To deal with that issue, Wagner et al. [2012] have proposed dedicated solutions, which jointly optimize the variables $\boldsymbol{\alpha}$ and \mathbf{e} , along with a deformation variable whose purpose is to automatically “realign images”.

Among other variants, it is worth mentioning some work using dictio-

nary learning to build the matrices \mathbf{D}_j [see, *e.g.*, Zhang and Li, 2010], and also the use of random projection features [Wright et al., 2009].

4.4 Patch classification and edge detection

The capability of dictionary learning for modeling particular types of signals has been exploited for classification tasks in several ways. Specifically, there are two main successful strategies; the first one learns one dictionary per class of signal and compare reconstruction errors obtained by the dictionaries in the same vein as the SRC classifier of the previous section does. The other direction uses sparse codes produced by a single dictionary as a high-level representation of signals for a subsequent classification procedure. We successively present these two lines of research, along with some successful applications.

Finding the dictionary that best represents a class of signals. Let us consider a training set of signals $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_k]$ in $\mathbb{R}^{m \times n}$ where the columns of each sub-matrix \mathbf{X}_j in $\mathbb{R}^{m \times n_j}$ correspond to signals of the j -th class. Even though dictionary learning is an unsupervised learning technique, it can be used to model each class independently by learning a dictionary \mathbf{D}_j “adapted” to the data \mathbf{X}_j . For instance, it is possible to define \mathbf{D}_j as the output of a learning algorithm for optimizing

$$\min_{\mathbf{D}_j \in \mathcal{C}, \mathbf{A}_j \in \mathbb{R}^{p \times n_j}} \frac{1}{2} \|\mathbf{X}_j - \mathbf{D}_j \mathbf{A}_j\|_F^2 + \lambda \sum_{i=1}^{n_j} \|\alpha_i\|_1,$$

where $\mathbf{A}_j = [\alpha_1, \dots, \alpha_{n_j}]$. Then, a new test signal \mathbf{x} in \mathbb{R}^m is classified according to the rule

$$\hat{j} = \arg \min_{j \in \{1, \dots, k\}} \left[\min_{\alpha \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}_j \alpha\|_2^2 + \lambda \|\alpha\|_1 \right]. \quad (4.7)$$

This simple rule was proposed by Ramirez et al. [2010], following earlier work of Mairal et al. [2008a] that focused on the ℓ_0 -penalty and discriminative formulations. Even though the dictionaries are learned in an unsupervised manner, (4.7) performs surprisingly well for data that is well modeled by dictionary learning. For instance, the approach

of Ramirez et al. [2010] essentially builds upon (4.7) while also encouraging the dictionaries \mathbf{D}_j to be mutually incoherent, and obtains competitive results for digit recognition, namely about 1.2% error rate on the MNIST dataset. In contrast, a nonlinear support vector machine (SVM) with a Gaussian kernel achieves 1.4% and a K -nearest neighbor approach 3.9% [LeCun et al., 1998a]. Of course, better results have been obtained by other techniques on the MNIST dataset, but they are in general not based on generic classifiers. More precisely, state-of-the-art algorithms typically encode invariance to image deformations [Ranzato et al., 2007, Bruna and Mallat, 2013].

Since dictionary learning is supposed to be well adapted to natural image patches, the classification rule (4.7) can naturally be applied to model the local appearance of particular object classes. For instance, we show in Figure 4.2 some results obtained with a formulation that is closely related to (4.7) for a weakly supervised learning task [see Mairal et al., 2008a, for more details]. Because the local appearance of bicycles is highly discriminative, the method performs well at detecting patches that overlap with bicycles in the test images.



Figure 4.2: Results of pixelwise classification of test images from Mairal et al. [2008a]. Two dictionaries are learned beforehand, one from images containing a bicycle, and one from images corresponding to background only. For each of the four test images displayed here, we also display a confidence map for the presence of a bicycle. For each pixel, a confidence value is computed for a patch centered at the pixel, by comparing the residual of the reconstruction errors obtained with the two dictionaries; yellow bright pixel values represent high confidence, whereas dark red pixel values represent low confidence.

Another successful application of dictionary learning is edge detection. A first attempt by Mairal et al. [2008d] is based on the principles described in this paragraph, and consists of learning two dictionaries, one on patches centered on edges, and one centered on background. The method was trained on the manually annotated Berkeley segmentation dataset [Martin et al., 2001] and was shown to perform as well as the dedicated approach called “Pb” [Martin et al., 2004]. Later, state-of-the-art results have been obtained by Xiaofeng and Bo [2012] by essentially combining two main ideas: (i) learning a single dictionary and exploiting the sparse codes for classification, a paradigm that will be the focus of the next paragraph, (ii) a dedicated way of performing multiscale feature pooling, giving more stability and some invariance to the sparse representations. We refer to Xiaofeng and Bo [2012] for all details and we present some results obtained with their software package in Figure 4.3.⁵ In the next paragraph, we now move to the second paradigm, namely exploiting the sparse codes as a nonlinear transformation of the input data for a subsequent linear classifier.



Figure 4.3: Edge detection results obtained with the software package of Xiaofeng and Bo [2012] on the Berkeley segmentation dataset BSD500 [Martin et al., 2001]. Four pairs of images are presented. The left ones represent input images and the right ones confidence maps for the presence of edges. Dark values represent high confidence. Best seen by zooming on a computer screen.

Using sparse codes for classification. Given a signal \mathbf{x} in \mathbb{R}^m , the sparse coding principle produces a sparse vector $\boldsymbol{\alpha}^*$ such that $\mathbf{x} \approx \mathbf{D}\boldsymbol{\alpha}^*$

⁵The software package is available here: <http://homes.cs.washington.edu/~xren>.

for some dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$. The process transforming \mathbf{x} into $\boldsymbol{\alpha}^*$ is highly nonlinear, and the entries of the vector $\boldsymbol{\alpha}^*$ can be interpreted as the “contributions” of the corresponding columns of \mathbf{D} to the reconstruction of \mathbf{x} .

When the dictionary atoms represent discriminative features for a classification task, it becomes appealing to use the vectors $\boldsymbol{\alpha}^*$ as new inputs to a classifier, *e.g.*, linear SVM or logistic regression. The strategy is thus significantly different than the one presented in the previous paragraph. The proof of concept was first demonstrated by Huang and Aviyente [2006] with a pre-defined fixed dictionary and a cost function inspired from linear discriminant analysis [see Hastie et al., 2009]. Later, dictionary learning was successfully used by Raina et al. [2007] for unsupervised feature learning as an effective way of exploiting unlabeled data for various classification tasks with image and text modalities.

Then, joint cost functions that involve both a discriminative term and a classical dictionary learning formulation have been proposed several times in the literature [Rodriguez and Sapiro, 2008, Mairal et al., 2008b, Pham and Venkatesh, 2008, Yang et al., 2011]. Let us consider a training set $(\mathbf{x}_i, y_i)_{i=1}^n$, where the vectors \mathbf{x}_i in \mathbb{R}^m are input signals and the scalars y_i are associated labels. The simplest cost functions that combine sparse coding and classification have the following form

$$\min_{\substack{\mathbf{D} \in \mathcal{C}, \mathbf{W} \in \mathbb{R}^{p \times k} \\ \mathbf{A} \in \mathbb{R}^{p \times n}}} \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1 + \gamma L(y_i, \mathbf{W}, \boldsymbol{\alpha}_i) \right] + \frac{\mu}{2} \|\mathbf{W}\|_{\mathbb{F}}^2, \quad (4.8)$$

where the loss function L measures how far a prediction based on $\boldsymbol{\alpha}_i$ is from the correct label y_i . The matrix \mathbf{W} represents model parameters for the classifier and the quadratic term $(\mu/2)\|\mathbf{W}\|_{\mathbb{F}}^2$ is a regularizer. Typically, L is a convex function and the prediction model is linear. For example, when the labels y_i are in $\{-1, +1\}$ for a binary classification problem, we may want to learn a linear decision rule parameterized by a vector \mathbf{w} (meaning the matrix \mathbf{W} has a single column), and use one

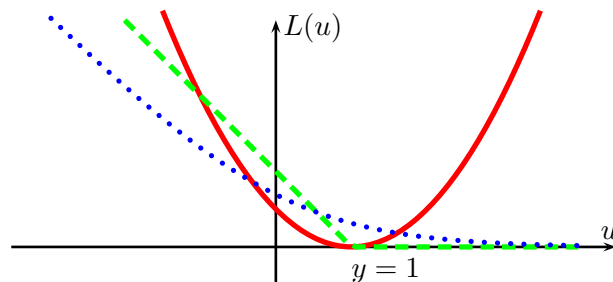


Figure 4.4: Visualization of three loss functions for a positive label $y = 1$. The plain red curve represents the square loss $L(u) = (1/2)(y - u)^2$, the dotted blue one represents the logistic loss $L(u) = \log(1 - e^{-yu})$, and the green dotted one represents the hinge loss $L(u) = \max(0, 1 - yu)$.

of the following loss functions

$$L(y, \mathbf{w}, \boldsymbol{\alpha}) \triangleq \begin{cases} \log(1 + e^{-y\mathbf{w}^\top \boldsymbol{\alpha}}) & \text{(logistic loss),} \\ \max(0, 1 - y\mathbf{w}^\top \boldsymbol{\alpha}) & \text{(hinge loss),} \\ \frac{1}{2}(y - \mathbf{w}^\top \boldsymbol{\alpha})^2 & \text{(square loss).} \end{cases} \quad (4.9)$$

For simplicity, we omit the fact that the loss function may involve an intercept—that is a constant term b in \mathbb{R} such that the linear model is $y \approx \text{sign}(\mathbf{w}^\top \boldsymbol{\alpha} + b)$ instead of simply $y \approx \text{sign}(\mathbf{w}^\top \boldsymbol{\alpha})$. Depending on the chosen loss function L , the formulation (4.8) can combine dictionary learning with a linear SVM, or a logistic regression classifier. We present the three loss functions (4.9) in Figure 4.4. The logistic and hinge loss are very similar; both of them are asymptotically linear and encourage the signs of the products $\mathbf{w}^\top \boldsymbol{\alpha}_i$ to be the same as y_i . In all cases, it is possible to iteratively decrease the value of the cost function (4.8) and obtain a stationary point by alternate minimization between the variables \mathbf{D} , \mathbf{A} , \mathbf{W} . As for the classical dictionary learning formulation, it is in fact impossible to find the global optimum; the problem is indeed convex with respect to each variable when keeping the other variables fixed, but it is not jointly convex.

The motivation of the formulation (4.8) is to encourage the dictionary to provide sparse decomposition vectors $\boldsymbol{\alpha}$ that are good for discrimination. It does so by encouraging them to have a small classification error on the training set measured by a convex loss function

from machine learning. There are however two possible caveats with such an approach, which we discuss now.

The sparse decomposition patterns may be unstable. The first caveat is the lack of stability of the estimator provided by sparse decomposition algorithms. For instance, the solution of the Lasso can significantly change after a small perturbation of the input data—that is, the vector

$$\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D}) \triangleq \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1, \quad (4.10)$$

assuming the optimization problem has a unique solution, can drastically change for small variations of \mathbf{x} , or even small variations of λ in pathological cases [Meinshausen and Bühlmann, 2010, Mairal and Yu, 2012]. For classification tasks, it is unfortunately often desirable to have stable representations, and thus improving this aspect may be important [Bruna and Mallat, 2013]. Several strategies have been proposed to deal with such an issue:

- **promoting incoherent dictionaries:** when $\mathbf{D}^\top \mathbf{D}$ is close to the identity matrix, the dictionary is called “incoherent”. When it is the case, sparse reconstruction algorithms may benefit from support recovery guarantees (see Section 1.5). Following this insight, Ramirez et al. [2009] have shown that promoting incoherent dictionaries may lead to better generalization properties. Later, this idea was also found useful in other classification schemes [Bo et al., 2013].
- **pooling features:** when a set of signals that are close to each other are available, *e.g.*, a set patches that significantly overlap with each other, feature pooling on the set of corresponding sparse codes can lead to more stable, but less sparse representations. For instance, Ren and Ramanan [2013], Bo et al. [2013] adopts such a strategy in an effective image representation for object detection.
- **replacing ℓ_1 by the elastic-net:** when replacing the ℓ_1 -norm in (4.10) by the Elastic-net penalty $\|\boldsymbol{\alpha}\|_1 + (\gamma/2)\|\boldsymbol{\alpha}\|_2^2$, it is possible to show that $\boldsymbol{\alpha}^*$ is Lipschitz continuous [Mairal et al., 2012],

with a Lipschitz constant proportional to $(1/\gamma)^2$. Here, the parameter γ controls a trade-off between sparsity and stability. As a side effect, it also leads to a unique solution of the sparse reconstruction problem due to the strong convexity of the penalty.

The label y is not available at test time. Interestingly, the formulation (4.8) is related to the super-resolution approach described in Section 3.4. even though the relation is only clear from hindsight. Indeed, the goal of classification is to map an input signal \mathbf{x} to an output label y . In contrast, super-resolution can be seen as a multivariate regression problem, whose goal is to learn a mapping between \mathbf{x} and a high-resolution signal—say, a signal \mathbf{y} . It is then not surprising that the classification formulation (4.8) is very similar to the regression one (3.6). In both cases, learning the mapping involves a dictionary \mathbf{D} for representing the input signals \mathbf{x} and a model that maps sparse codes $\boldsymbol{\alpha}$ to the variable to predict. In (4.8), the parameters \mathbf{W} play for instance the same role as the parameters \mathbf{D}_h in (3.6).

In Section 3.4, we have discussed some limitations of the super-resolution formulation (3.6), which are in fact also relevant for the classification task. Given some fixed \mathbf{D} and \mathbf{W} , we remark that the sparse codes $\boldsymbol{\alpha}_i$ for the training data may be obtained as follows

$$\boldsymbol{\alpha}_i \in \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 + \gamma L(y_i, \mathbf{W}, \boldsymbol{\alpha}). \quad (4.11)$$

Unfortunately, given a new test signal \mathbf{x} , the label y is not available and the formulation (4.10) has to be used instead to obtain the corresponding sparse code $\boldsymbol{\alpha}$. Therefore, there is a discrepancy in the way the signals are encoded during the training and test phases.

In Section 3.4, we have addressed that issue by presenting an extension involving a bilevel optimization problem. It is also possible to follow the same strategy for the classification task, but we develop this method in the next section since it draws a strong connection between dictionary learning and neural networks.

4.5 Connections with neural networks

There are obvious connections between dictionary learning and neural networks. Some of them have already appeared in this monograph with multilayer schemes involving feature maps, which are typical architectures of convolutional neural networks and related work [LeCun et al., 1998a, Serre et al., 2005]. In this section, we present two other important connections. First, we introduce the concept of “backpropagation” for dictionary learning, which is directly borrowed from neural networks [LeCun et al., 1998b], and which addresses an issue that was left opened at the end of the previous section. Then, we focus on particular networks whose purpose is to approximate efficiently the solution of sparse coding problems [Gregor and LeCun, 2010].

Backpropagation rules for dictionary learning. In Figure 4.5, we present two paradigms for dictionary learning coupled with a prediction task, *e.g.*, regression or classification. In the first one, the dictionary is obtained without supervision and the prediction model is learned in a subsequent step. In the latter paradigm, all parameters, including the dictionary, are learned jointly for the prediction task. This requires a similar concept as “backpropagation” [LeCun et al., 1998b], which consists of using the chain rule in neural networks.

Formally, let us consider a training set of signals $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$ associated to measurements $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$ in $\mathbb{R}^{k \times n}$ representing the variables to predict at test time. Such a setting is quite general on purpose: the vectors \mathbf{y}_i may represent for instance the high-resolution patches from Section 3.4 or the half-toned patches from Section 3.5, but they may also represent a label for the classification tasks of Section 4.4. In this last case, $k = 1$ and the \mathbf{y}_i ’s are simply scalars.

In the formulation of interest, the signals \mathbf{x}_i are encoded with \mathbf{D} :

$$\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D}) = \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 + \frac{\mu}{2} \|\boldsymbol{\alpha}\|_2^2, \quad (4.12)$$

where the use of the elastic-net penalty instead of ℓ_1 will be made clear in the sequel. Then, we consider a similar prediction model as in the previous section; following the same methodology, we introduce a loss

function L , and the corresponding empirical risk minimization problem

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{W} \in \mathbb{R}^{p \times k}} \frac{1}{n} \sum_{i=1}^n L(\mathbf{y}_i, \mathbf{W}, \boldsymbol{\alpha}^*(\mathbf{x}_i, \mathbf{D})) + \frac{\gamma}{2} \|\mathbf{W}\|_{\mathbb{F}}^2, \quad (4.13)$$

Specifically, any loss from the previous section is considered to be appropriate here. The advantage of (4.12) over (4.11) is that the encoding scheme remains the same between training and testing, and yet \mathbf{D} and \mathbf{W} can be learned jointly for the final prediction task in (4.13). Moreover (4.13) extends the super-resolution formulation of (3.8) to more general problems such as classification.

Even though (4.13) seems appealing, it is unfortunately particularly difficult to solve. In addition to being nonconvex, the relation between \mathbf{D} and the objective function goes through the argmin of a nonsmooth optimization problem. As a consequence, the cost function is also nonsmooth with respect to \mathbf{D} and it is not clear how to obtain a descent direction.

In an asymptotic regime where enough training data is available, Mairal et al. [2012] show that the cost function becomes differentiable, and they derive a stochastic gradient descent algorithm from this theoretical result. They call this approach “task-driven dictionary learning”. More precisely, consider the expected cost

$$f(\mathbf{D}, \mathbf{W}) \triangleq \mathbb{E}_{(\mathbf{y}, \mathbf{x})} [L(\mathbf{y}, \mathbf{W}, \boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D}))],$$

where the expectation is taken to the unknown probability distribution of the data pairs (\mathbf{y}, \mathbf{x}) . Under mild assumptions on the data distribution and when $\mu > 0$, it is possible to show that f is differentiable and that its gradient admits a closed form:

$$\begin{cases} \nabla_{\mathbf{W}} f(\mathbf{D}, \mathbf{W}) &= \mathbb{E}_{(\mathbf{y}, \mathbf{x})} \left[\nabla L(\mathbf{y}, \mathbf{W}, \boldsymbol{\alpha}^*)^\top \right], \\ \nabla_{\mathbf{D}} f(\mathbf{D}, \mathbf{W}) &= \mathbb{E}_{(\mathbf{y}, \mathbf{x})} \left[-\mathbf{D} \boldsymbol{\beta}^* \boldsymbol{\alpha}^{*\top} + (\mathbf{x} - \mathbf{D} \boldsymbol{\alpha}^*) \boldsymbol{\beta}^{*\top} \right], \end{cases} \quad (4.14)$$

where $\boldsymbol{\alpha}^*$ is short for $\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D})$,

$$\boldsymbol{\beta}^*[\Gamma^c] = 0, \text{ and } \boldsymbol{\beta}^*[\Gamma] \triangleq \left(\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma + \mu \mathbf{I} \right)^{-1} \nabla_{\boldsymbol{\alpha}[\Gamma]} L(\mathbf{y}, \mathbf{W}, \boldsymbol{\alpha}^*),$$

where Γ is the set of nonzero entries of $\boldsymbol{\alpha}^*$. Here, the role of parameter μ is to ensure that the matrix to be inverted is well conditioned. Since the gradient has the form of an expectation, it is natural to use a stochastic gradient descent algorithm to optimize the

cost function [Bottou and Bousquet, 2008]. Even though the nonconvexity makes it impossible to find the global optimum, it was shown by Mairal et al. [2012] that following heuristics from the neural network literature [LeCun et al., 1998b], this strategy yields sufficiently good results for many prediction tasks.

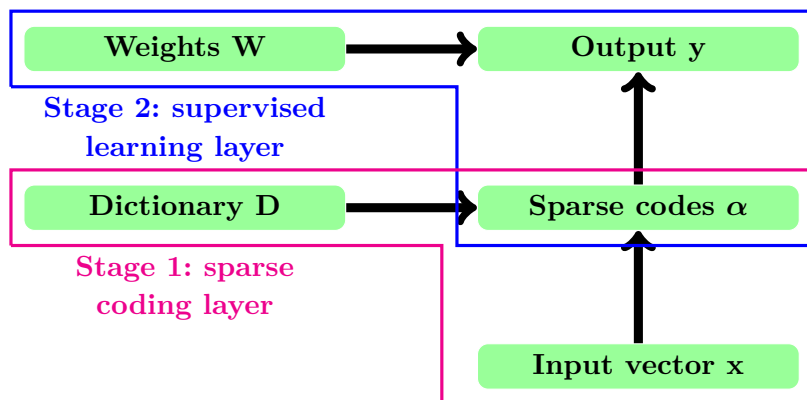
The concept of learning \mathbf{D} to be good for prediction is similar to the backpropagation principle in neural networks where all the parameters of the network are learned for a prediction task using the chain rule [LeCun et al., 1998b]. To the best of our knowledge, Bradley and Bagnell [2008] were the first to use such a terminology for dictionary learning; they proposed indeed a different supervised dictionary learning formulation with smoothed approximations of sparsity-inducing penalties. In a heuristic fashion, similar rules as (4.14) have also been derived by Boureau et al. [2010] and Yang et al. [2010b] in the context of classification and by Yang et al. [2012a] for super-resolution (see Section 3.4).

Fast approximations of sparse coding. Another link between neural networks and dictionary learning has been established by Kavukcuoglu et al. [2010b] and Gregor and LeCun [2010], who have trained neural networks to approximate sparse codes. Given some fixed dictionary \mathbf{D} and a training set \mathbf{X} , their idea is to learn a simple nonlinear function $g(\mathbf{W}, \mathbf{x})$ that approximates $\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D})$ for any signal \mathbf{x} of interest. By using similar notation as in the previous paragraphs, the general formulation proposed by Gregor and LeCun [2010] is the following

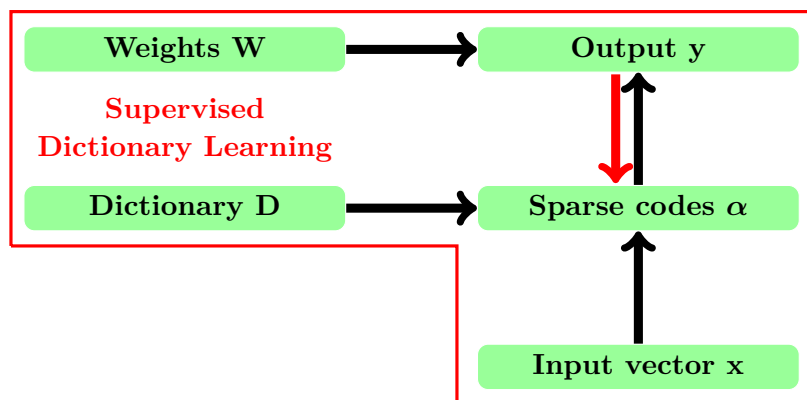
$$\min_{\mathbf{W} \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^n \|\boldsymbol{\alpha}^*(\mathbf{x}_i, \mathbf{D}) - g(\mathbf{W}, \mathbf{x}_i)\|_2^2,$$

where \mathbf{W} is a set of parameters for the prediction function g optimized with stochastic gradient descent and \mathcal{W} is an optimization domain. The basic prediction function g introduced by Kavukcuoglu et al. [2010b] is the composition of a linear transformation and a pointwise nonlinearity, which can be interpreted as a one-layer neural network:

$$g(\mathbf{W}, \mathbf{x})[j] \triangleq \gamma_j \tanh(\mathbf{w}_j^\top \mathbf{x} + b_j), \quad (4.15)$$



(a) Without backpropagation.



(b) With backpropagation.

Figure 4.5: Two paradigms for dictionary learning are illustrated in the figure. An input signal \mathbf{x} is sparsely encoded by using a dictionary \mathbf{D} , producing sparse codes α . Then, a second model with parameters \mathbf{W} makes a prediction \mathbf{y} based on α . In (a), the dictionary \mathbf{D} is learned unsupervised in a first layer—that is, without exploiting output information \mathbf{y} . Then, the parameters \mathbf{W} of the prediction model are learned afterwards in the second layer. In (b), we illustrate the concept of supervised dictionary learning, where \mathbf{D} and \mathbf{W} are jointly learned for the prediction task; this requires “backpropagating” information, represented by the red arrow, from the top layer to the bottom one.

where the scalars γ_j and b_j are some weights that can be also optimized. The motivation for such an approach is to reduce the cost of computing sparse codes at test time. Whereas obtaining $\boldsymbol{\alpha}^*(\mathbf{x}, \mathbf{D})$ requires solving an optimization problem that can be costly (see Section 5), the complexity of computing $g(\mathbf{W}, \mathbf{w})$ is the same as a matrix-vector multiplication, making it possible to develop real-time applications.

However, there are two major shortcomings to the original approach of Kavukcuoglu et al. [2010b], as noted by Gregor and LeCun [2010]. First, the choice of the non-linearity \tanh does not yield exactly sparse vectors, but it can be replaced by more appropriate shrinkage functions. Second, the scheme (4.15) encodes each entry of the output vector $g(\mathbf{W}, \mathbf{x})$ independently and cannot possibly take into account the fact that there exists correlation among the dictionary elements. To address this issue, Gregor and LeCun [2010] have proposed an iterative approach with a higher complexity than (4.15), but with a better approximation (see their paper for more details).

4.6 Other applications

Because the research topic has been very active in the last decade, sparse estimation and dictionary learning have been used in numerous ways in computer vision; unfortunately, writing an exhaustive review would be an endless exercise, and we have instead focused so far on a few approaches related to visual recognition in images. In the remaining paragraphs, we aim to be slightly more exhaustive and briefly mention a few other successful applications.

Action recognition in videos. We start with natural extensions to videos of visual recognition pipelines originally developed for images. Castrodad and Sapiro [2012] and Guha and Ward [2012] have proposed two related approaches for action classification in movies. We present here the main principles from the method of Castrodad and Sapiro [2012], which has shown competitive results in several standard evaluation benchmarks.

Let us consider n training video sequences represented by vectors

$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ in \mathbb{R}^{lT} , where T is the number of frames in the sequence and l is the number of pixels in each frame. Each sequence contains a subject performing a particular action and the goal is to predict the type of action in subsequent test videos. We assume that there are k different possible actions.

To address this classification problem, Castrodad and Sapiro [2012] have introduced a two-layer dictionary learning scheme, after pre-processing the videos to replace the raw pixels by temporal gradients. Then, the two layers of the pipeline can be described as follows:

- **local model of actions:** it is common in videos to work with three-dimensional patches with a time dimension to describe the local appearance of actions [Laptev, 2005]. For each class j , one dictionary \mathbf{D}_j is learned on the training data to represent the appearance of such patches within the action class j . This step is particularly useful when some actions are locally discriminative.
- **encoding of full sequences:** once the dictionaries \mathbf{D}_j are learned, a train or test video is processed by extracting its overlapping three-dimensional patches, and by encoding them with the concatenated dictionary $\mathbf{D} = [\mathbf{D}_1, \dots, \mathbf{D}_k]$, say with p dictionary elements. This produces a representation α_i for every patch i of the sequence. These patches are pooled into a single vector β in \mathbb{R}^p for the sequence.
- **global model of actions:** the first layer provides a high-dimensional vector for each video. The second layer is exploiting these nonlinear representations of the sequences as inputs to the classification scheme (4.7). As a result, the layer is learning a global model for each action class, by learning again one dictionary per class. The difference with the first layer is that the dictionaries model the appearance of full sequences instead of local three-dimensional patches.

Note that, to simplify, we have also omitted some useful heuristics such as the removing of three-dimensional patches with low energy, corresponding to motionless regions. We have also not provided all mathematical details such as the chosen dictionary learning formulations

and algorithms. We refer to Castrodad and Sapiro [2012] for a more exhaustive description of the method.

Visual tracking. We now move to a technique for visual tracking introduced by Mei and Ling [2009], which relies on the sparsity-inducing effect of the ℓ_1 -norm and on some ideas from the robust face recognition system of Wright et al. [2009]. Given a stream of images $\mathbf{x}_1, \mathbf{x}_2, \dots$, the goal of tracking is to find a bounding box around a moving object in the sequence. Typically, the bounding box is annotated for the first image, and the task is to find the boxes in the subsequent frames by exploiting temporally consistency.

The approach of Mei and Ling [2009] is inspired by template matching techniques [Lucas and Kanade, 1981, Matthews et al., 2004]. It assumes that at time t , a set of p templates $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ in $\mathbb{R}^{m \times p}$ is available, where each column \mathbf{d}_j is a small image filled by the object of interest—in other words, a bounding box. Typically, each \mathbf{d}_j has been either obtained by a ground truth annotation at the beginning of the sequence, or is one of the estimated boxes in the previous frames. Then, estimating the position of the object in frame t consists of scanning all possible bounding boxes in that frame, and comparing how well they can be reconstructed by the templates \mathbf{D} .

Concretely, let us assume that there are l possible boxes, and let us denote by \mathbf{R}_k the linear operator that extracts the k -th box from the image and subsamples it to make a representation with m parameters. Then, the optimization problem is simply

$$\hat{k} \in \arg \min_{k \in \{1, \dots, l\}} \left[\min_{\alpha \in \mathbb{R}_+^p, \mathbf{e} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{R}_k \mathbf{x}_t - \mathbf{D} \alpha - \mathbf{e}\|_2^2 + \lambda \|\mathbf{e}\|_1 + \sum_{j=1}^p \eta_j |\alpha[j]| \right], \quad (4.16)$$

where \mathbf{e} models the occlusion as in the face recognition formulation of Wright et al. [2009] presented in Section 4.3. The ℓ_1 -norm encourages \mathbf{e} to be sparse, whereas the weighted- ℓ_1 -norm also promotes the selection of a few templates from \mathbf{D} . The role of the weights η_j is to give more importance to some templates, especially the ones that have been selected recently in the previous frames.

Overall, the tracking scheme performs two successive operations

for each frame. It selects the new bounding box with (4.16), and then updates the dictionary \mathbf{D} and the weights η_j , according to an ad-hoc procedure [see Mei and Ling, 2009, for the details]. According to a recent survey [Wu et al., 2013], trackers based on sparse representations, which are essentially improved versions of Mei and Ling [2009], perform well compared to the state of the art, and are able to deal with occlusion effectively.

Data visualization. Finally, we deviates from visual recognition by presenting applications of sparse estimation to data visualization in computer vision [Elhamifar et al., 2012, Chen et al., 2014]. The first approach proposed by Elhamifar et al. [2012] can be viewed as an extension of the sparse subspace clustering method of Elhamifar and Vidal [2009], and was successfully applied to the problem of video summarization.

More precisely, consider a video sequence $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T]$ in $\mathbb{R}^{l \times T}$ where each \mathbf{x}_t in \mathbb{R}^l is the t -th frame. The goal of video summarization is to find a set of $s \ll T$ frames that best “represents” the video. The selection problem is formulated as follows

$$\min_{\mathbf{A} \in \mathbb{R}^{T \times T}} \|\mathbf{X} - \mathbf{X}\mathbf{A}\|_{\text{F}}^2 \quad \text{s.t.} \quad \Psi_0(\mathbf{A}) \leq s \quad \text{and} \quad \forall i, \sum_{j=1}^T \alpha_i[j] = 1, \quad (4.17)$$

where $\mathbf{A} = [\alpha_1, \dots, \alpha_T]$ and $\Psi_0(\mathbf{A})$ counts the number of rows of \mathbf{A} that are non-zero. The constraints have two simultaneous effects: first, they encourage every frame \mathbf{x}_t to be close to a linear combination $\mathbf{X}\alpha_t$ of other frames with coefficients α_t that sum to one; second, since \mathbf{A} has at most s non-zero rows, only s frames from \mathbf{X} can be used in the approximations $\mathbf{X}\alpha_t$. The video sequence is consequently “summarized” by these s frames.

To visualize large image collections, Chen et al. [2014] has used related ideas by exploiting an older unsupervised learning technique called archetypal analysis, which is presented in Section 2.6. Specifically, let us consider a collection of images, denoted again by $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, but which do not necessarily form a video sequence. Archetypal analysis looks for a dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$ that “represents”

well the data,

$$\min_{\substack{\alpha_i \in \Delta_p \text{ for } 1 \leq i \leq n \\ \beta_j \in \Delta_n \text{ for } 1 \leq j \leq p}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2 \quad \text{s.t.} \quad \forall j, \mathbf{d}_j = \mathbf{X}\beta_j,$$

where Δ_p and Δ_n are simplex sets—that is, sets of vectors with non-negative entries that sum to one. In contrast to the original dictionary learning formulation, the dictionary elements, called archetypes, are convex combinations of data points and admit a simple interpretation.

In Figure 4.6, we present a result obtained by Chen et al. [2014] on a database of images downloaded on the Flickr website with the request “Paris”. Each image is encoded using the Fisher vector representation [Perronnin et al., 2010] and $p = 256$ archetypes are learned. Each archetype \mathbf{d}_j is a convex combination $\mathbf{X}\beta_j$ of a few input images with a sparse vector β_j due to the simplicial constraint. For each archetype, it is natural to visualize images that are used in the sparse combination, as well as the corresponding non-zero coefficients $\beta_j[i]$. To some extent, the resulting figure “summarizes” the image collection.

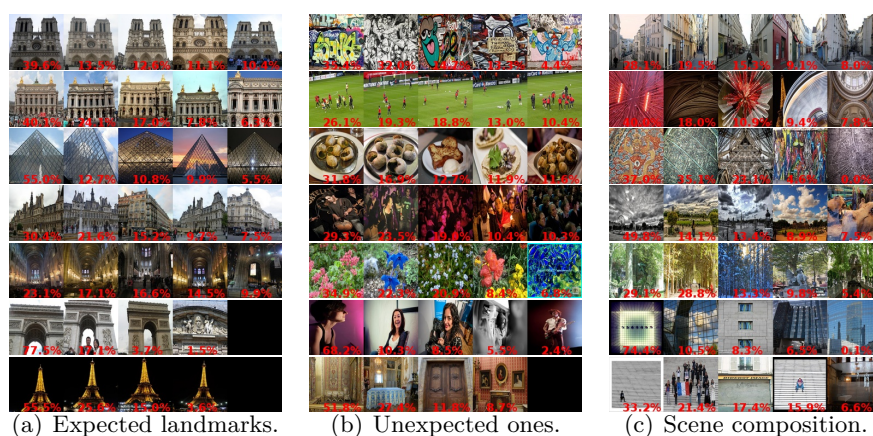


Figure 4.6: A few archetypes learned from 36 600 pictures corresponding to the request “Paris” downloaded from Flickr and sorted by “relevance”. The figure is produced by Chen et al. [2014] with the software package SPAMS available at <http://spams-devel.gforge.inria.fr/> The numbers in red represent the values of the non-zero coefficients $\beta_j[i]$ for each image. In (a), we obtain expected landmarks such as the Eiffel tower or Notre Dame; in (b), we display some unexpected archetypes corresponding to soccer, graffiti, food, flowers and social gatherings; in (c), we see some archetypes that do not seem to have some semantic meaning, but that capture some scene compositions or texture present in the dataset. Best seen by zooming on a computer screen.

5

Optimization Algorithms

In the previous parts of the monograph, many formulations for dealing with image processing and computer vision tasks have been introduced, leading in general to cost functions to optimize. However, we have not presented yet in great details how the corresponding algorithms should be implemented in practice. The purpose of this section is to introduce basic optimization tools to address sparse estimation and dictionary learning problems. Because the literature is vast [see Bach et al., 2012a, for a review], we focus on a few techniques that we consider to be easy to use and efficient enough.

Specifically, we mostly present algorithms that are parameter-free. This means that in addition to *model parameters* that are inherent to the formulation, they do not require tuning any input value to converge; they should at most require choosing a number of iterations or a stopping criterion. We also put the emphasis on algorithms that are easy to implement: a few lines of a high-level language, *e.g.*, Matlab or Python, are in general sufficient to obtain an effective prototype, even though cutting-edge implementations may be more involved to develop—that is, they may require a low-level language such as C/C++ and optimized

third-party libraries for linear algebra such as BLAS and LAPACK.¹

First, we address classical ℓ_0 -regularized problems, before treating the case of the ℓ_1 -penalty. Then, we present iterative reweighted ℓ_1 and ℓ_2 techniques, and conclude with dictionary learning algorithms. Many of these methods have been used for the experiments conducted in this monograph, and are available in the SPAMS toolbox.²

5.1 Sparse reconstruction with the ℓ_0 -penalty

We start our review of optimization techniques with least square problems that are regularized with the ℓ_0 -penalty. More precisely, we are interested in finding approximate solutions of

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_0 \leq k, \quad (5.1)$$

or

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\boldsymbol{\alpha}\|_0 \quad \text{s.t.} \quad \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \leq \varepsilon, \quad (5.2)$$

where \mathbf{x} is an input signal in \mathbb{R}^m and \mathbf{D} is a dictionary in $\mathbb{R}^{m \times p}$. There are mainly two classes of algorithms for dealing with such NP-hard problems: greedy algorithms and iterative hard-thresholding methods. We start with algorithms of the first kind, which are designed for dealing with both formulations. They are commonly known as *matching pursuit* and *orthogonal matching pursuit* in signal processing, and *forward selection* techniques in statistics.

Coordinate descent algorithm: matching pursuit (MP). Introduced by Mallat and Zhang [1993], matching pursuit is a coordinate descent algorithm that iteratively selects one entry of the current estimate $\boldsymbol{\alpha}$ at a time and updates it. It is closely related to projection pursuit algorithms from the statistics literature [Friedman and Stuetzle, 1981], and is presented in Algorithm 2, where $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_p]$ is a dictionary with unit-norm columns. Each iteration performs in fact a one dimensional minimization of the residual with respect to the \hat{j} -th entry of $\boldsymbol{\alpha}$

¹see <http://www.netlib.org/blas/>.

²<http://spams-devel.gforge.inria.fr/>.

Algorithm 2 Matching pursuit algorithm.

Require: Signal \mathbf{x} in \mathbb{R}^m , dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$ with unit-norm columns, and stopping criterion k or ε .

- 1: Initialize $\boldsymbol{\alpha} \leftarrow 0$;
- 2: **while** $\|\boldsymbol{\alpha}\|_0 < k$ if k is provided or $\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 > \varepsilon$ if ε is provided
do
- 3: select the coordinate with maximum partial derivative:

$$\hat{j} \in \arg \max_{j=1 \dots p} |\mathbf{d}_j^\top (\mathbf{x} - \mathbf{D}\boldsymbol{\alpha})|;$$

- 4: update the coordinate

$$\boldsymbol{\alpha}[\hat{j}] \leftarrow \boldsymbol{\alpha}[\hat{j}] + \mathbf{d}_j^\top (\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}); \quad (5.3)$$

- 5: **end while**

- 6: **return** the sparse decomposition $\boldsymbol{\alpha}$ in \mathbb{R}^p .
-

when keeping the other entries fixed. The update (5.3) can indeed be interpreted as

$$\boldsymbol{\alpha}[\hat{j}] \leftarrow \arg \min_{\alpha \in \mathbb{R}} \left\| \mathbf{x} - \sum_{l \neq \hat{j}} \boldsymbol{\alpha}[l] \mathbf{d}_l - \alpha \mathbf{d}_{\hat{j}} \right\|_2^2. \quad (5.4)$$

Therefore, the residual monotonically decreases during the optimization, whereas the model sparsity $\|\boldsymbol{\alpha}\|_0$ increases or remains constant until the stopping criterion is satisfied. One drawback of matching pursuit is that the same entry \hat{j} can be selected several times during the optimization, possibly leading to a large number of iterations. The next algorithm, called “orthogonal matching pursuit”, does not suffer from this issue, and usually provides a better sparse approximation with additional computational cost.

Active-set algorithm: orthogonal matching pursuit (OMP). The algorithm is similar in spirit to matching pursuit, but enforces the residual to be always orthogonal to all previously selected variables, which is equivalent to saying that the algorithm reoptimizes the value of all

non-zero coefficients once it selects a new variable. As a result, each iteration of OMP is more costly than those of MP, but the procedure is guaranteed to converge in a finite number of iterations. The algorithm appears in the signal processing literature [Pati et al., 1993], and is called “forward selection” in statistics (see Section 1.1).

Several variants of OMP have been proposed that essentially differ in the way a new variable is selected at each iteration. In this monograph, we focus on an effective variant introduced by Gharavi-Alkhansari and Huang [1998], Cotter et al. [1999], which is called “order recursive orthogonal matching pursuit”. To simplify, we simply use the terminology “OMP” to denote this variant and omit the term “order recursive” in the rest of the monograph. The procedure is presented in Algorithm 3, where \mathbf{D}_Γ represents the matrix of size $\mathbb{R}^{m \times |\Gamma|}$ whose columns are those of \mathbf{D} indexed by Γ , and Γ^c denotes the complement set of Γ in $\{1, \dots, p\}$. The classical OMP algorithm (not the order recursive variant) chooses \hat{j} exactly as in the MP algorithm. Through our experience, we have found that the order recursive variant provides a slightly better sparse approximation in practice, even though it may seem more complicated to implement at first sight.

Unlike MP, the set Γ of non-zero variables of $\boldsymbol{\alpha}$ strictly increases by one unit after each iteration, and it is easy to show that the residual is always orthogonal to the matrix \mathbf{D}_Γ of previously selected dictionary elements. Indeed, we have at each step of the algorithm:

$$\boldsymbol{\alpha}_\Gamma = \left(\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma \right)^{-1} \mathbf{D}_\Gamma^\top \mathbf{x},$$

(assuming \mathbf{D}_Γ to be full-rank), and thus

$$\mathbf{D}_\Gamma^\top (\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}) = \mathbf{D}_\Gamma^\top (\mathbf{x} - \mathbf{D}_\Gamma \boldsymbol{\alpha}_\Gamma) = 0.$$

It is interesting to note that efficient implementations of Algorithm 3 exist, despite the fact that the algorithm naively requires to minimize $|\Gamma^c|$ quadratic functions—equivalently solve $|\Gamma^c|$ linear systems—at each iteration. The main ingredients for a fast implementation have been discussed by Cotter et al. [1999], and are summarized below:

1. if the Cholesky factorization of $(\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma)$ is already computed at the beginning of an iteration, finding the index \hat{j} to select in

Algorithm 3 Orthogonal matching pursuit - order recursive variant.

Require: Signal \mathbf{x} in \mathbb{R}^m , dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$, and stopping criterion k or ε .

- 1: Initialize the active set $\Gamma = \emptyset$ and $\boldsymbol{\alpha} \leftarrow 0$;
- 2: **while** $|\Gamma| < k$ if k is provided or $\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 > \varepsilon$ if ε is provided
do
- 3: select a new coordinate \hat{j} that leads to the smallest residual:

$$(\hat{j}, \hat{\boldsymbol{\beta}}) \in \arg \min_{j \in \Gamma^c, \boldsymbol{\beta} \in \mathbb{R}^{|\Gamma|+1}} \left\| \mathbf{x} - \mathbf{D}_{\Gamma \cup \{j\}} \boldsymbol{\beta} \right\|_2^2; \quad (5.5)$$

- 4: update the active set and the solution $\boldsymbol{\alpha}$:

$$\begin{aligned} \Gamma &\leftarrow \Gamma \cup \{\hat{j}\}; \\ \boldsymbol{\alpha}_\Gamma &\leftarrow \hat{\boldsymbol{\beta}} \quad \text{and} \quad \boldsymbol{\alpha}_{\Gamma^c} \leftarrow 0; \end{aligned}$$

- 5: **end while**

- 6: **return** the sparse decomposition $\boldsymbol{\alpha}$ in \mathbb{R}^p .
-

the update (5.5) can be found with very few operations by using simple linear algebra relations [Cotter et al., 1999].

2. given some index \hat{j} , we have

$$\hat{\boldsymbol{\beta}} = (\mathbf{D}'_\Gamma{}^\top \mathbf{D}'_\Gamma)^{-1} \mathbf{D}'_\Gamma{}^\top \mathbf{x},$$

where $\Gamma' = \Gamma \cup \{\hat{j}\}$. Whereas the matrix inversion naively costs $O(|\Gamma'|^3)$ operations, it can be obtained in $O(|\Gamma|^2)$ operations if the Cholesky factorization $(\mathbf{D}'_\Gamma{}^\top \mathbf{D}'_\Gamma)^{-1}$ is already computed.

3. if $\mathbf{Q} \triangleq \mathbf{D}^\top \mathbf{D}$ is precomputed, a significant gain in terms of speed can be obtained by exploiting the fact that the quadratic function $\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2$ expands as $\|\mathbf{x}\|_2^2 - 2\mathbf{q}^\top \boldsymbol{\alpha} + \boldsymbol{\alpha}^\top \mathbf{Q}\boldsymbol{\alpha}$. This trivial reformulation leads to an implementation that never explicitly computes the residual $\mathbf{r} = \mathbf{x} - \mathbf{D}\boldsymbol{\alpha}$, but instead works with the quantity $\mathbf{D}^\top \mathbf{r}$, initialized with \mathbf{q} , and updated at every iteration. This simple strategy may be useful when large set of signals $\mathbf{x}_1, \dots, \mathbf{x}_n$ needs to be encoded in parallel with the same

dictionary. Then, it is worth computing \mathbf{Q} once for all, before processing all signals. It is then possible to obtain an implementation of Algorithm 3 whose complexity for encoding these n signals is

$$\underbrace{O(mp^2)}_{\text{compute } \mathbf{D}^\top \mathbf{D} \text{ once}} + \underbrace{nO(k^3)}_{\text{Cholesky factorizations}} + \underbrace{nO(pm + pk^2)}_{\text{update of } \mathbf{D}^\top \mathbf{r}}.$$

A similar strategy for the classical OMP algorithm that involves a simpler selection rule of the variable \hat{j} is discussed by Rubinstein et al. [2008].

Even though the algorithm only provides an approximate solution of the ℓ_0 -regularized least square problem in general, OMP was shown by Tropp and Gilbert [2007] to be able to reliably recover a sparse signal from random measurements. Interestingly, the conditions that are required on the dictionary are similar to the ones for sparse recovery in the compressed sensing literature (see Section 1.5).

Gradient-descent technique: iterative hard-thresholding. A significantly different approach than the previous greedy algorithms consists of using a gradient descent principle, an approach called *iterative hard-thresholding* [Starck et al., 2003, Herrity et al., 2006, Blumensath and Davies, 2009], designed to find an approximate solution of either (5.1) or

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_0, \quad (5.6)$$

for some penalty parameter λ . The method is presented in Algorithm 4 and is an instance of the more general *proximal gradient descent* algorithms [see Bach et al., 2012a, for a review].

As we will see, when η is smaller than the inverse of the largest eigenvalue of $\mathbf{D}^\top \mathbf{D}$, the algorithm monotonically decreases the value of the objective function. To see this, we can interpret the iterative hard-thresholding algorithm as a majorization-minimization scheme [Lange et al., 2000], which consists of minimizing at each iteration a locally tight upper-bound of the objective. This principle is illustrated in Figure 5.1.

Algorithm 4 Iterative hard thresholding.

Require: Signal \mathbf{x} in \mathbb{R}^m , dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$, target sparsity k or penalty parameter λ , step size η , number of iterations T , initial solution α_0 in \mathbb{R}^p (with $\|\alpha_0\|_0 \leq k$ if k is provided).

- 1: Initialize $\alpha \leftarrow \alpha_0$;
- 2: **for** $t = 1, \dots, T$ **do**
- 3: perform one step of gradient descent:

$$\alpha \leftarrow \alpha + \eta \mathbf{D}^\top (\mathbf{x} - \mathbf{D}\alpha);$$

- 4: choose threshold τ to be the k -th largest magnitude among the entries of α if k is provided, or $\tau = \sqrt{2\lambda}$ if λ is provided;
- 5: **for** $j = 1, \dots, p$ **do**
- 6: hard-thresholding:

$$\alpha[j] \leftarrow \begin{cases} \alpha[j] & \text{if } |\alpha[j]| \geq \tau, \\ 0 & \text{otherwise.} \end{cases}$$

- 7: **end for**
 - 8: **end for**
 - 9: **return** the sparse decomposition α in \mathbb{R}^p .
-

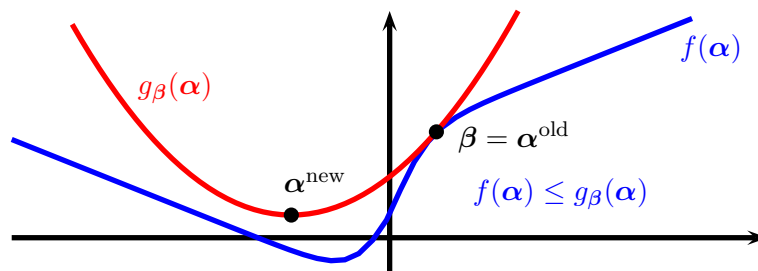


Figure 5.1: Illustration of the basic majorization-minimization principle. The objective function f (in blue) needs to be minimized. At the current iteration, a majorizing surrogate g_β (in red) is computed that is locally tight at the previous estimate α^{old} . After minimizing the surrogate, we obtain a new estimate α^{new} such that $f(\alpha^{\text{new}}) \leq f(\alpha^{\text{old}})$.

The majorizing surrogate can be obtained from the following inequality that holds for all $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ in \mathbb{R}^p :

$$\begin{aligned}
\frac{1}{2}\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 &= \frac{1}{2}\|\mathbf{x} - \mathbf{D}\boldsymbol{\beta} - \mathbf{D}(\boldsymbol{\alpha} - \boldsymbol{\beta})\|_2^2 \\
&= \frac{1}{2}\|\mathbf{x} - \mathbf{D}\boldsymbol{\beta}\|_2^2 - (\boldsymbol{\alpha} - \boldsymbol{\beta})^\top \mathbf{D}^\top (\mathbf{x} - \mathbf{D}\boldsymbol{\beta}) + \frac{1}{2}\|\mathbf{D}(\boldsymbol{\alpha} - \boldsymbol{\beta})\|_2^2 \\
&\leq \frac{1}{2}\|\mathbf{x} - \mathbf{D}\boldsymbol{\beta}\|_2^2 - (\boldsymbol{\alpha} - \boldsymbol{\beta})^\top \mathbf{D}^\top (\mathbf{x} - \mathbf{D}\boldsymbol{\beta}) + \frac{1}{2\eta}\|\boldsymbol{\alpha} - \boldsymbol{\beta}\|_2^2 \\
&= \underbrace{C_\beta + \frac{1}{2\eta}\left\|\boldsymbol{\alpha} - \left(\boldsymbol{\beta} + \eta\mathbf{D}^\top(\mathbf{x} - \mathbf{D}\boldsymbol{\beta})\right)\right\|_2^2}_{g_\beta(\boldsymbol{\alpha})},
\end{aligned} \tag{5.7}$$

where C_β is a term independent of $\boldsymbol{\alpha}$. All equalities are simple linear algebra relations and the inequality comes from the fact that $\mathbf{z}^\top \mathbf{D}^\top \mathbf{D} \mathbf{z} \leq \frac{1}{\eta} \mathbf{z}^\top \mathbf{z}$ for all \mathbf{z} in \mathbb{R}^p , since $1/\eta$ is assumed to be larger than the largest eigenvalue of $\mathbf{D}^\top \mathbf{D}$. Under this assumption, we have the relation $\frac{1}{2}\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \leq g_\beta(\boldsymbol{\alpha})$ for all $\boldsymbol{\alpha}$ in \mathbb{R}^p , with an equality for $\boldsymbol{\alpha} = \boldsymbol{\beta}$. Then, we have two possible cases:

- if the problem of interest is (5.1), k is provided to Algorithm 4, and it is easy to show that the value of $\boldsymbol{\alpha}$ at the end of one iteration is the solution of

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} g_\beta(\boldsymbol{\alpha}) \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_0 \leq k,$$

where $\boldsymbol{\beta}$ is the value of $\boldsymbol{\alpha}$ at the beginning of the iteration.

- if one is interested in the penalized problem (5.6), the value of $\boldsymbol{\alpha}$ after each iteration of Algorithm 4 is instead the solution of

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} g_\beta(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|_0,$$

and the function $\boldsymbol{\alpha} \mapsto g_\beta(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|_0$ is a majorizing surrogate of the objective $\frac{1}{2}\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_0$.

In both cases, Algorithm 4 can be interpreted as a majorization-minimization scheme and the value of the objective function monotonically decreases.

5.2 Sparse reconstruction with the ℓ_1 -norm

In the previous section, we have studied three optimization problems (5.1), (5.2) and (5.6). Similarly, we consider three formulations of interest involving the ℓ_1 -norm:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1, \quad (5.8)$$

and

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\boldsymbol{\alpha}\|_1 \quad \text{s.t.} \quad \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \leq \varepsilon, \quad (5.9)$$

and

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_1 \leq \mu. \quad (5.10)$$

We have also presented three optimization strategies to deal with the ℓ_0 -penalty, namely: matching pursuit, orthogonal matching pursuit, and iterative hard-thresholding algorithms. Interestingly, it is possible to deal with the ℓ_1 -penalty by following similar techniques, even though the links between ℓ_1 -approaches and ℓ_0 -ones are only clear from hindsight.

Coordinate descent. The counterpart of matching pursuit for ℓ_1 -regularized problems is probably the coordinate descent method. Even though the general scheme is classical in optimization [see Bertsekas, 1999], it was first proposed by Fu [1998] for minimizing the Lasso (5.8). Its convergence to the exact solution of the convex non-smooth problem such as (5.8) was shown by Tseng and Yun [2009], and its empirical efficiency was demonstrated by Wu and Lange [2008]. The coordinate descent approach is presented in Algorithm 5, where S_λ denotes the soft-thresholding operator already introduced in Section 1.2:

$$S_\lambda : \theta \mapsto \text{sign}(\theta) \max(|\theta| - \lambda, 0),$$

and the selection rule for choosing \hat{j} can be on the following strategies:

- cycle in $\{1, \dots, p\}$ (the most typical choice);
- pick up \hat{j} uniformly at random in $\{1, \dots, p\}$;

Algorithm 5 Coordinate descent algorithm for the Lasso (5.8).

Require: Signal \mathbf{x} in \mathbb{R}^m , dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$, penalty parameter λ , number of iterations T , initial solution $\boldsymbol{\alpha}_0$ in \mathbb{R}^p .

- 1: Initialize $\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha}_0$;
- 2: **for** $t = 1, \dots, T$ **do**
- 3: select one coordinate \hat{j} , *e.g.*, by cycling or at random;
- 4: update the coordinate by soft-thresholding:

$$\boldsymbol{\alpha}[\hat{j}] \leftarrow S_\lambda \left(\boldsymbol{\alpha}[\hat{j}] + \frac{\mathbf{d}_j^\top (\mathbf{x} - \mathbf{D}\boldsymbol{\alpha})}{\|\mathbf{d}_j\|_2^2} \right); \quad (5.11)$$

- 5: **end for**
 - 6: **return** the sparse decomposition $\boldsymbol{\alpha}$ in \mathbb{R}^p .
-

- choose the entry \hat{j} exactly as in the matching pursuit algorithm.

It is relatively easy to show that the update (5.11) is in fact minimizing the objective (5.8) with respect to the \hat{j} -th entry of $\boldsymbol{\alpha}$ when keeping the others fixed. More precisely, (5.11) is equivalent to

$$\boldsymbol{\alpha}[\hat{j}] \leftarrow \arg \min_{\alpha \in \mathbb{R}} \left\| \mathbf{x} - \sum_{l \neq \hat{j}} \boldsymbol{\alpha}[l] \mathbf{d}_l - \alpha \mathbf{d}_j \right\|_2^2 + \lambda |\alpha|,$$

which is very similar to the matching pursuit update (5.4) but with the ℓ_1 -penalty in the formulation. The benefits of the coordinate descent scheme are three-fold: (i) it is extremely simple to implement; (ii) it has no parameter; (iii) it can be very efficient in practice [see Bach et al., 2012a, Section 8, for a benchmark].

Assuming that one chooses the cycling selection rule, the complexity of a naive update of $\boldsymbol{\alpha}[j]$ is $O(mp)$ for computing the matrix-vector multiplication $\mathbf{D}\boldsymbol{\alpha}$. Fortunately, better complexities can be obtained by adopting one of the two following alternative strategies:

1. during the algorithm, we choose to update an auxiliary variable representing the residual \mathbf{r} such that we always have $\mathbf{r} = \mathbf{x} - \mathbf{D}\boldsymbol{\alpha}$. Given \mathbf{r} , updating $\boldsymbol{\alpha}[\hat{j}]$ only cost $O(m)$ operations because of the inner-product $\mathbf{d}_j^\top \mathbf{r}$. Then, updating \mathbf{r} either has zero cost

if $\alpha[\hat{j}]$ does not change, or cost $O(m)$ otherwise. Finally, the per-iteration complexity of this strategy is always $O(m)$;

2. assuming that the matrix $\mathbf{D}^\top \mathbf{D}$ is pre-computed, it is possible to obtain an even better per-iteration complexity by updating an auxiliary variable \mathbf{z} instead of \mathbf{r} such that we always have $\mathbf{z} = \mathbf{D}^\top (\mathbf{x} - \mathbf{D}\alpha)$. Given \mathbf{z} , updating $\alpha[j]$ only cost $O(1)$ operations. Updating \mathbf{z} has either zero cost if $\alpha[j]$ does not change, or cost $O(p)$ operations. As a result, the per-iteration complexity is at most $O(p)$, but more interestingly, each time a variable $\alpha[j]$ is equal to zero before its update and remains equal to zero after the update (a typical scenario for very sparse solutions), the cost per iteration is $O(1)$.

Extensions of coordinate descent to deal with more general smooth loss functions and group-Lasso penalties have also been proposed [Tseng and Yun, 2009]. These variants are also easy to implement, and exhibit good performance in practice [see Bach et al., 2012a].

Proximal gradient algorithm - iterative soft-thresholding. The second approach that we consider is the ℓ_1 -counterpart of the iterative hard-thresholding algorithm called “iterative soft-thresholding” [Nowak and Figueiredo, 2001, Figueiredo and Nowak, 2003, Starck et al., 2003, Daubechies et al., 2004]. Originally designed to address (5.8), it can be easily modified to deal with (5.10). It is similar to the hard-thresholding technique of the previous section and is presented in Algorithm 6, where η satisfies the same property as in the previous section—that is, $1/\eta$ is larger than the largest eigenvalue of $\mathbf{D}^\top \mathbf{D}$. Note that the orthogonal projection (5.12) onto the ℓ_1 -ball can be done efficiently in linear time [Brucker, 1984, Duchi et al., 2008].

Then, the method can also be interpreted under the majorization-minimization point of view illustrated in Figure 5.1. Indeed, let us consider the two cases:

- if the problem of interest is (5.8), the user provides the parameter λ , and it is easy to show that the value of α at the end of one

Algorithm 6 Proximal gradient algorithm for (5.8) or (5.10).

Require: Signal \mathbf{x} in \mathbb{R}^m , dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$, regularization parameter λ or μ , gradient descent step size η , number of iterations T , initial solution $\boldsymbol{\alpha}_0$ in \mathbb{R}^p (with $\|\boldsymbol{\alpha}_0\|_1 \leq T$ if T is provided).

- 1: Initialize $\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha}_0$;
- 2: **for** $t = 1, \dots, T$ **do**
- 3: perform one step of gradient descent:

$$\boldsymbol{\alpha} \leftarrow \boldsymbol{\alpha} + \eta \mathbf{D}^\top (\mathbf{x} - \mathbf{D}\boldsymbol{\alpha});$$

- 4: **if** λ is provided **then**
- 5: soft-thresholding: for all $j = 1, \dots, p$,

$$\boldsymbol{\alpha}[j] \leftarrow S_\lambda(\boldsymbol{\alpha}[j]);$$

- 6: **else**
- 7: μ is provided and one performs the orthogonal projection

$$\boldsymbol{\alpha} \leftarrow \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left[\frac{1}{2} \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|_2^2 \text{ s.t. } \|\boldsymbol{\beta}\|_1 \leq \mu \right]; \quad (5.12)$$

- 8: **end if**
 - 9: **end for**
 - 10: **return** the sparse decomposition $\boldsymbol{\alpha}$ in \mathbb{R}^p .
-

iteration is the solution of

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} g_{\beta}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|_1,$$

where g_{β} is defined in (5.7), and the function $\boldsymbol{\alpha} \mapsto g_{\beta}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|_1$ is a majorizing surrogate of the objective $(1/2)\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1$.

- if one is interested instead in the penalized problem (5.10), the value of $\boldsymbol{\alpha}$ after each iteration of Algorithm 4 is the solution of

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} g_{\beta}(\boldsymbol{\alpha}) \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_1 \leq \mu,$$

where β is the value of $\boldsymbol{\alpha}$ at the beginning of the iteration.

The majorization-minimization principle explains why the algorithm monotonically decreases the value of the objective function, but a more refined analysis can provide convergence rates. In general, we know that the objective function value converges to the optimal one with the convergence rate $O(1/t)$, and that some variants such as FISTA admit a better rate $O(1/t^2)$, as shown by Beck and Teboulle [2009], Nesterov [2013]. These methods are relatively easy to implement, and choosing η is not an issue in practice. When an appropriate value for η is not available beforehand, it is possible to automatically find a good one with a line-search scheme.

Similarly as coordinate descent, extensions to more general problems than ℓ_1 -regularized least squares are possible. In general proximal gradient methods are indeed adapted to solving convex problems of the form

$$\min_{\boldsymbol{\alpha} \in \mathcal{A}} f(\boldsymbol{\alpha}) + \psi(\boldsymbol{\alpha}),$$

where \mathcal{A} is a convex set in \mathbb{R}^p , f is differentiable and its gradient is uniformly Lipschitz continuous, and ψ is a convex nonsmooth function such that the so-called *proximal operator* of ψ

$$\arg \min_{\boldsymbol{\alpha} \in \mathcal{A}} \frac{1}{2} \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|_2^2 + \psi(\boldsymbol{\alpha}),$$

can be efficiently computed for all $\boldsymbol{\beta}$ in \mathbb{R}^p . As a result, proximal gradient methods have been successfully used with the Group Lasso penalty including the hierarchical and structured variants of Section 1.3 [see Jenatton et al., 2011b, Mairal et al., 2010b].

Homotopy and LARS algorithms Finally, the homotopy method [Osborne et al., 2000, Efron et al., 2004] is related to orthogonal matching pursuit, but for solving the ℓ_1 -regularized problems (5.8), (5.9), and (5.10). Similarly to OMP, the algorithm maintains an active-set of variables Γ , initialized with the empty set, and iteratively updates it by one variable at a time.

The homotopy method follows in fact the regularization path of the Lasso—that is, it finds the set of all solutions of (5.8) for all values of the regularization parameter λ , formally defined as

$$\mathcal{P} \triangleq \{\boldsymbol{\alpha}^*(\lambda) : \lambda > 0\},$$

where $\boldsymbol{\alpha}^*(\lambda)$ is the solution of (5.8), which is assumed to be unique under a small technical assumption [see Osborne et al., 2000, Mairal and Yu, 2012]. A remarkable property of the path is that it can be shown to be piecewise linear. In other words, it can be described by a finite number k of linear segments, as illustrated earlier in Figure 1.4:

$$\mathcal{P} = \bigcup_{l=1}^k \{\theta \boldsymbol{\alpha}^*(\lambda_l) + (1 - \theta) \boldsymbol{\alpha}^*(\lambda_{l+1}) : \theta \in [0, 1]\}. \quad (5.13)$$

The homotopy method finds a point on the path, computes the direction of the current linear segment, and follows it until the direction of the path changes. The algorithm either returns the full regularization path—that is, the sequence of $\boldsymbol{\alpha}^*(\lambda_l)$ from Eq. (5.13), or stops if it reaches a solution for a target λ , a target residual ϵ , or ℓ_1 -norm μ . Thus, it can be used for solving the three variants (5.8), (5.9), or (5.10).

The piecewise linearity of the path was discovered by Markowitz [1952] in the context of portfolio selection, and was turned into an effective algorithm for solving the Lasso by Osborne et al. [2000] and Efron et al. [2004]. It should be also noted that solving (5.8) for all values of λ is in fact an instance of *parametric quadratic programming*, for which path following algorithm encompassing the homotopy method appear early in the optimization literature [Ritter, 1962]. The homotopy turns out to be extremely efficient for solving very sparse medium-scale problems [see Bach et al., 2012a], even though its worst-case complexity has been shown to be exponential [Mairal and Yu,

2012]. Similar to the simplex algorithm for solving linear programs, the empirical complexity is therefore extremely bad in the worst-case scenario, but it is most often very good in practical ones.

Following the presentation of Mairal and Yu [2012], we describe the method in Algorithm 7. The formula for the direction of the path is given in Eq. (5.14), where the quantity $\boldsymbol{\eta}$ is defined as $\boldsymbol{\eta} \triangleq \text{sign}(\mathbf{D}^\top(\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}^*))$ in $\{-1, 0, 1\}^p$. This relation can be derived from optimality conditions of the Lasso presented in Section 1.3, but the derivation requires technical developments that are beyond the scope of the monograph. We refer instead to Mairal and Yu [2012] for more details. Interestingly, implementing efficiently Algorithm 7 raises similar difficulties as for implementing orthogonal matching pursuit. We need indeed to be able to update the inverse matrices $(\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma)^{-1}$ at each iteration when the active-set Γ changes. The Cholesky decomposition or the Woodbury formula [see Golub and Van Loan, 2012] can be used for that purpose.

For example, let us give a few details on how the Woodbury formula can be useful. Denoting by $\Gamma' \triangleq \Gamma \cup \{\hat{j}\}$, we have the classical relation

$$\left(\mathbf{D}_{\Gamma'}^\top \mathbf{D}_{\Gamma'}\right)^{-1} = \begin{bmatrix} \left(\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma\right)^{-1} + \frac{1}{s} \mathbf{z} \mathbf{z}^\top & -\frac{1}{s} \mathbf{z} \\ -\frac{1}{s} \mathbf{z}^\top & \frac{1}{s} \end{bmatrix},$$

where $s \triangleq \mathbf{d}_j^\top \mathbf{d}_j - \mathbf{d}_j^\top (\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma)^{-1} \mathbf{d}_j$ is called the *Schur complement*, and $\mathbf{z} \triangleq (\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma)^{-1} \mathbf{d}_j$. Thus, computing $(\mathbf{D}_{\Gamma'}^\top \mathbf{D}_{\Gamma'})^{-1}$ given $(\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma)^{-1}$ from the previous iteration can be done by a few matrix-vector multiplications only. In the same vein, obtaining $(\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma)^{-1}$ from $(\mathbf{D}_{\Gamma'}^\top \mathbf{D}_{\Gamma'})^{-1}$ is also a matter of manipulating a few linear algebra relations.

We also remark that the implementation of the homotopy can benefit from the pre-computation of the matrix $\mathbf{D}^\top \mathbf{D}$ in the exact same way as for orthogonal matching pursuit, such that both implementations have in fact the same complexity (up to a constant factor that is slightly larger for the homotopy method).

5.3 Iterative reweighted- ℓ_1 methods

We have mentioned in the introduction of this monograph that non-convex sparsity-inducing penalties are often used in sparse estimation.

Algorithm 7 Homotopy algorithm for the Lasso.

Require: Signal \mathbf{x} in \mathbb{R}^m , dictionary \mathbf{D} in $\mathbb{R}^{m \times p}$.

- 1: $\lambda \leftarrow \|\mathbf{D}^\top \mathbf{x}\|_\infty$;
- 2: $\Gamma \leftarrow \{j \in \{1, \dots, p\} : |\mathbf{d}_j^\top \mathbf{x}| = \lambda\}$;
- 3: **while** $\lambda > 0$ or according to another stopping criterion **do**
- 4: compute the current direction of the regularization path:

$$\boldsymbol{\alpha}_\Gamma^*(\lambda) = (\mathbf{D}_\Gamma^\top \mathbf{D}_\Gamma)^{-1} (\mathbf{D}_\Gamma^\top \mathbf{x} - \lambda \boldsymbol{\eta}_\Gamma) \quad \text{and} \quad \boldsymbol{\alpha}_{\Gamma^c}^*(\lambda) = 0. \quad (5.14)$$

- 5: find the smallest $\tau > 0$ such that one of the following events occurs:
 - there exists j in Γ^c such that $|\mathbf{d}_j^\top (\mathbf{x} - \mathbf{D} \boldsymbol{\alpha}^*(\lambda - \tau))| = \lambda - \tau$;
then, add j to Γ ;
 - there exists j in Γ such that $\boldsymbol{\alpha}^*(\lambda - \tau)[j]$ hits zero; then,
remove j from Γ .

It is also easy to show that the value of τ can be obtained in closed form such that one can “jump” from a kink to another;

- 6: update $\lambda \leftarrow \lambda - \tau$; record the pair $(\lambda, \boldsymbol{\alpha}^*(\lambda))$;
 - 7: **end while**
 - 8: **Return:** sequence of recorded values $(\lambda, \boldsymbol{\alpha}^*(\lambda))$.
-

These regularization functions typically have the form

$$\psi(\boldsymbol{\alpha}) = \sum_{j=1}^p \varphi(|\boldsymbol{\alpha}[j]|), \quad (5.15)$$

where the functions φ are concave, non-decreasing and differentiable on \mathbb{R}_+ . A natural way of dealing with such penalties when they appear in the objective function is to use an approach called DC-programming, where DC stands for “difference of convex” [see Gasso et al., 2009, for a review]. Even though the resulting optimization problems are non-convex and impossible to solve exactly in a reasonable amount of time, DC-programming is a simple technique to obtain a stationary point by using the majorization-minimization principle already illustrated in Figure 5.1. It consists of exploiting a majorizing surrogate of ψ obtained by linearizing the function φ in (5.15), which is upper-bounded by its linear approximations because of concavity, as shown in Figure 5.2.

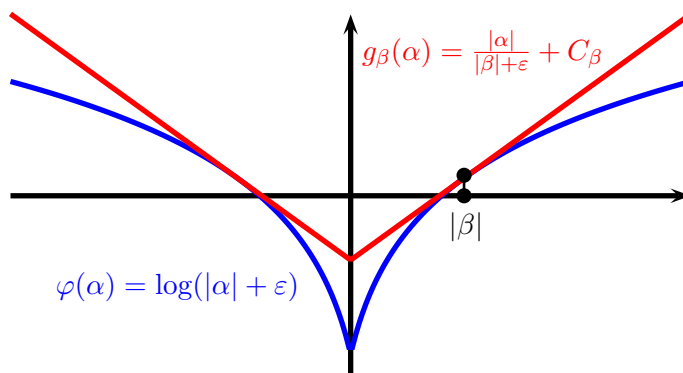


Figure 5.2: Illustration of the DC-programming approach. The function φ (in blue) is upper-bounded by its linear approximation (in red), which is tight at the point β . The resulting majorizing surrogate is a weighted ℓ_1 -norm plus a constant C_β .

For example, for minimizing

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \sum_{j=1}^p \log(|\boldsymbol{\alpha}[j]| + \varepsilon),$$

the resulting majorization-minimization algorithm solves a sequence of

weighted Lasso problems at each iteration

$$\boldsymbol{\alpha}^{\text{new}} \leftarrow \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \sum_{j=1}^p \frac{|\boldsymbol{\alpha}[j]|}{|\boldsymbol{\alpha}^{\text{old}}[j]| + \varepsilon} \right], \quad (5.16)$$

where $\boldsymbol{\alpha}^{\text{old}}$ denotes the current estimate at the beginning of the iteration. The update (5.16) is a reweighted- ℓ_1 algorithm, which is known to provide sparser solutions than the regular Lasso. Such majorization-minimization approaches for non-convex sparse estimation have been introduced in various contexts. For example, Fazel [2002], Fazel et al. [2003] and then later Candès et al. [2008] use such a principle for the regularization function $\varphi : x \mapsto \log(x + \varepsilon)$, whereas Figueiredo and Nowak [2005], Figueiredo et al. [2007] propose majorization-minimization algorithms for dealing with the ℓ_q -penalty with $q < 1$.

5.4 Iterative reweighted- ℓ_2 methods

We now focus on another type of algorithm for dealing with convex regularizers such as the ℓ_1 -norm. Similar to the reweighted- ℓ_1 -algorithms of the previous section, the approaches we present here consist of solving a sequence of sub-problems, each one involving a simple penalty. The key idea is that many non-smooth convex regularizers can be seen as minima of quadratic functions, leading to reweighted ℓ_2 -algorithms.

Variational formulation of the ℓ_1 -norm. We start with the simplest case of the ℓ_1 -norm, which admits the following variational form [Grandvalet and Canu, 1999, Daubechies et al., 2010]:

$$\|\boldsymbol{\alpha}\|_1 = \inf_{\boldsymbol{\eta} \in \mathbb{R}_+^p} \frac{1}{2} \sum_{j=1}^p \left\{ \frac{\boldsymbol{\alpha}[j]^2}{\boldsymbol{\eta}[j]} + \boldsymbol{\eta}[j] \right\},$$

with the optimal $\boldsymbol{\eta} \in \mathbb{R}_+^p$ obtained as $\boldsymbol{\eta}[j] = |\boldsymbol{\alpha}[j]|$. The Lasso optimization problem then becomes:

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \min_{\boldsymbol{\eta} \in \mathbb{R}_+^p} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \frac{\lambda}{2} \sum_{j=1}^p \left\{ \frac{\boldsymbol{\alpha}[j]^2}{\boldsymbol{\eta}[j]} + \boldsymbol{\eta}[j] \right\}.$$

Fortunately, this problem is jointly convex in $(\boldsymbol{\alpha}, \boldsymbol{\eta})$. The minimization with respect to $\boldsymbol{\eta}$ given $\boldsymbol{\alpha}$ can be done in closed form (by taking absolute values of the components of $\boldsymbol{\alpha}$); the minimization with respect to $\boldsymbol{\alpha}$ given $\boldsymbol{\eta}$ is a weighted least-squares problem, which can be solved by the solution of a linear system as:

$$\boldsymbol{\alpha} = [\mathbf{D}^\top \mathbf{D} + \lambda \text{Diag}(\boldsymbol{\eta})^{-1}]^{-1} \mathbf{D}^\top \mathbf{x}.$$

Note that since the linear system above is ill-conditioned when some entries of $\boldsymbol{\eta}$ are very small, it is better to equivalently reformulate it as

$$\boldsymbol{\alpha} = \text{Diag}(\boldsymbol{\eta})^{1/2} [\text{Diag}(\boldsymbol{\eta})^{1/2} \mathbf{D}^\top \mathbf{D} \text{Diag}(\boldsymbol{\eta})^{1/2} + \lambda \mathbf{I}]^{-1} \text{Diag}(\boldsymbol{\eta})^{1/2} \mathbf{D}^\top \mathbf{x}.$$

This naturally leads to alternating minimization procedures where one alternates between minimizing with respect to $\boldsymbol{\eta}$ as above in closed form, and minimizing with respect to $\boldsymbol{\alpha}$, which is now easier because the regularizer has become a quadratic function. With quadratic losses, this leads to a linear system for which many existing algorithms exist. However, because the function $(\boldsymbol{\alpha}[j], \boldsymbol{\eta}[j]) \mapsto \boldsymbol{\alpha}[j]^2 / \boldsymbol{\eta}[j]$ cannot be extended to be continuous at $(0, 0)$, alternating minimization does not converge and it is customary to add a penalty term $\sum_{j=1}^p \frac{\varepsilon}{2\boldsymbol{\eta}[j]}$ with $\varepsilon > 0$ small, which leads to an update for $\boldsymbol{\eta}[j]$ of the form $\boldsymbol{\eta}[j] = \sqrt{|\boldsymbol{\alpha}[j]|^2 + \varepsilon}$ that avoids instabilities. See more details in Bach et al. [2012a, Section 5].

Extensions to ℓ_q -quasi-norms. More generally, for any $q \in (0, 2)$, and $\|\boldsymbol{\alpha}\|_q = (\sum_{j=1}^q |\boldsymbol{\alpha}[j]|^q)^{1/q}$ and $r = q/(2 - q)$, we have [Jenatton et al., 2010b]:

$$\|\boldsymbol{\alpha}\|_q = \inf_{\boldsymbol{\eta} \in \mathbb{R}_+^p} \frac{1}{2} \sum_{j=1}^p \frac{\boldsymbol{\alpha}[j]^2}{\boldsymbol{\eta}[j]} + \frac{1}{2} \|\boldsymbol{\eta}\|_r,$$

with the optimal $\boldsymbol{\eta}$ in \mathbb{R}_+^p obtained as $\boldsymbol{\eta}[j] = |\boldsymbol{\alpha}[j]|^{2-q} \|\boldsymbol{\alpha}\|_q^{q-1}$. Note that similar variational formulations may be obtained for the squared quasi-norms [Micchelli and Pontil, 2005]. This gives an alternative to reweighted ℓ_1 -formulations presented in Section 5.3 for $q < 1$ (the non-convex case).

Extensions to group norms. The various formulations above can be extended to structured sparsity-inducing norms, such as the ones

presented in Section 1.3. For example, for the group-lasso norm in Eq. (1.20), we have:

$$\sum_{g \in \mathcal{G}} \|\alpha[g]\|_2 = \inf_{\eta \in \mathbb{R}_+^p} \frac{1}{2} \sum_{g \in \mathcal{G}} \left\{ \frac{\|\alpha[g]\|_2^2}{\eta[g]} + \eta[g] \right\},$$

which also leads to a quadratic function of α and to simple algorithms based on solving linear systems when the loss is quadratic. Note that these quadratic variational formulations extend to all norms and lead to multiple kernel learning formulations when used with positive-definite kernels [see more details in Bach et al., 2012a].

5.5 Optimization for dictionary learning

Finally, we review a few dictionary learning algorithms that leverage the optimization algorithms for sparse estimation presented in the previous sections. We start with the stochastic gradient descent method [see Kushner and Yin, 2003, Bottou and Bousquet, 2008, and references therein] since it is very close to the original algorithm of Olshausen and Field [1997]. Then, we move to other classical approaches such as alternate minimization [Engan et al., 1999, Lee et al., 2007], block coordinate descent, K-SVD [Aharon et al., 2006], and then present efficient methods based on stochastic approximations [Mairal et al., 2010a, Skretting and Engan, 2010].

We consider two dictionary learning formulations. For both of them, the goal is to learn a dictionary \mathbf{D} in \mathcal{C} , where \mathcal{C} is the set of matrices in $\mathbb{R}^{m \times p}$ whose columns are constrained to have less than unit ℓ_2 -norm. Then, we use a sparsity-inducing penalty ψ either as a penalty

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2 + \lambda \psi(\alpha_i), \quad (5.17)$$

or as a constraint

$$\min_{\mathbf{D} \in \mathcal{C}, \mathbf{A} \in \mathbb{R}^{p \times n}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2 \quad \text{s.t.} \quad \psi(\alpha_i) \leq \mu, \quad (5.18)$$

and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$ is the training set of signals. Even though these problems are non-convex, many algorithms have been developed

with empirical good performance for various tasks, *e.g.*, image denoising. In general, these methods have no other guarantee than providing a stationary point (in the best case), but they have been used successfully in practical contexts, as shown in other parts of this monograph.

Stochastic gradient descent. The first algorithm proposed by Olshausen and Field [1996, 1997] addresses the formulation (5.17), where ψ is the ℓ_1 -norm or a smooth approximate sparsity-inducing penalty such as $\psi(\boldsymbol{\alpha}) = \sum_{j=1}^p \log(\boldsymbol{\alpha}[j]^2 + \varepsilon)$. The algorithm can be described as a heuristic stochastic gradient descent, which alternates between two stages: (i) gradient steps with a fixed step size computed on mini-batches of the training set and (ii) rescaling heuristic that prevents the norm of the columns \mathbf{d}_j to grow out of bounds—thus, taking implicitly the constraint $\mathbf{D} \in \mathcal{C}$ into account.

A less heuristic but very related procedure is the projected stochastic gradient descent method, which is presented in Algorithm 8. This approach relies on the following equivalent formulation to (5.17):

$$\min_{\mathbf{D} \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{D}), \quad (5.19)$$

where

$$L(\mathbf{x}_i, \mathbf{D}) \triangleq \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1. \quad (5.20)$$

At each iteration, the algorithm selects one signal \hat{i} and performs one gradient step $\mathbf{D} \leftarrow \mathbf{D} - \eta_t \nabla_{\mathbf{D}} L(\mathbf{x}_{\hat{i}}, \mathbf{D})$ with a step size η_t . It can indeed be shown that the function $\mathbf{D} \rightarrow L(\mathbf{x}_{\hat{i}}, \mathbf{D})$ is differentiable according to Danskin's theorem [see Bertsekas, 1999, Proposition B.25] and that its gradient admits a closed form $\nabla_{\mathbf{D}} L(\mathbf{x}_{\hat{i}}, \mathbf{D}) = -(\mathbf{x}_{\hat{i}} - \mathbf{D}\boldsymbol{\alpha}_{\hat{i}})\boldsymbol{\alpha}_{\hat{i}}^\top$, where $\boldsymbol{\alpha}_{\hat{i}}$ is defined in (5.21). Then, the algorithm performs a Euclidean projection $\Pi_{\mathcal{C}}$ onto the convex set \mathcal{C} , which corresponds to the renormalization update $\mathbf{d}_j \rightarrow (1/\max(\|\mathbf{d}_j\|_2, 1))\mathbf{d}_j$ for all j in $\{1, \dots, p\}$. A practical variant also consists of using mini-batches instead of drawing a single data point at each iteration. For simplicity, we omit this variant in our description of Algorithm 8.

The stochastic gradient descent approach is effective [see Mairal et al., 2010a, for a benchmark], but it raises a few practical

difficulties such as choosing well the step sizes η_t in a data-independent manner. Classical strategies use for instance step sizes of the form $\eta_t = \eta/(t + t_0)^\gamma$, but finding an appropriate set of hyper-parameters η, t_0 , and γ may be challenging in practice.

Algorithm 8 Stochastic gradient descent for (5.17).

Require: Signals $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$, initial dictionary \mathbf{D}_0 in \mathcal{C} , penalty parameter λ , number of iterations T .

- 1: Initialize $\mathbf{D} \leftarrow \mathbf{D}_0$;
- 2: **for** $t = 1, \dots, T$ **do**
- 3: select one signal \mathbf{x}_i at random;
- 4: compute the sparse code:

$$\boldsymbol{\alpha}_i \in \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \psi(\boldsymbol{\alpha}); \quad (5.21)$$

- 5: perform one projected gradient descent step:

$$\mathbf{D} \leftarrow \Pi_{\mathcal{C}} \left[\mathbf{D} - \eta_t (\mathbf{D}\boldsymbol{\alpha}_i - \mathbf{x}_i) \boldsymbol{\alpha}_i^\top \right];$$

- 6: **end for**
 - 7: **return** the sparse decomposition $\boldsymbol{\alpha}$ in \mathbb{R}^p .
-

Alternate minimization. The most classical approach for dictionary learning is the alternate minimization scheme, which consists of optimizing (5.17) or (5.18) by alternating between two minimization steps: one with respect to \mathbf{D} with the sparse codes \mathbf{A} fixed, and one with respect to \mathbf{A} with \mathbf{D} fixed. In practice, this approach is not as fast as a well-tuned stochastic gradient descent algorithm, but it is parameter-free and we subjectively find it very reliable according to our experience for image processing and computer vision tasks.

Alternate minimization was first proposed by Engan et al. [1999] for dealing with the ℓ_0 -penalty with the method of optimal directions (MOD), and was revisited later by Lee et al. [2007] with the ℓ_1 -regularization. We present the MOD approach in Algorithm 9 and its ℓ_1 -counterpart in Algorithm 10. As we shall see, they slightly differ in

the dictionary update step, which can be simplified for the ℓ_0 -penalty.

The sparse codes α_i for MOD are obtained with any algorithm for dealing approximately with ℓ_0 , *e.g.*, any technique from Section 5.1. Since all the vectors α_i can be computed in parallel, it is worth mentioning that this sparse decomposition step can benefit from an efficient implementation of OMP including the heuristics of Section 5.1, notably those consisting of pre-computing $\mathbf{D}^\top \mathbf{D}$ before updating all the α_i 's in parallel. For dealing with ℓ_1 , the homotopy method can be used similarly, with the same heuristics that apply to OMP.

Algorithm 9 Method of optimal directions for $\psi = \ell_0$.

Require: Signals $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$, initial dictionary \mathbf{D}_0 in \mathcal{C} , regularization parameter μ , number of iterations T .

1: Initialize $\mathbf{D} \leftarrow \mathbf{D}_0$;

2: **for** $t = 1, \dots, T$ **do**

3: compute the sparse codes, *e.g.*, with OMP:

4: **for** $i = 1, \dots, n$ **do**

5:

$$\alpha_i \approx \arg \min_{\alpha \in \mathbb{R}^p} \left[\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\alpha\|_2^2 \text{ s.t. } \|\alpha\|_0 \leq \mu \right];$$

6: **end for**

7: update the dictionary \mathbf{D} :

$$\mathbf{D} \leftarrow \Pi_{\mathcal{C}}[\mathbf{X}\mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}]; \quad (5.22)$$

8: **end for**

9: **return** the dictionary \mathbf{D} in \mathcal{C} .

Moving on now to the dictionary update step, we remark that (5.22) and (5.23) are not the same and are not equivalent to each other, even though, at first sight, the optimization problem with respect to \mathbf{D} when \mathbf{A} is fixed seems to be independent of the regularization ψ . The MOD update indeed minimizes the least-square function (5.23) by discarding the constraint \mathbf{D} in \mathcal{C} , yielding a solution $\mathbf{X}\mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top)^{-1}$ (assuming $\mathbf{A}\mathbf{A}^\top$ to be invertible), before projecting it onto the constraint set \mathcal{C} . In constrained optimization [see Bertsekas, 1999], it is usually not

Algorithm 10 Alternate minimization for $\psi = \ell_1$.

Require: Signals $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$, initial dictionary \mathbf{D}_0 in \mathcal{C} , regularization parameter λ or μ , number of iterations T .

- 1: Initialize $\mathbf{D} \leftarrow \mathbf{D}_0$;
- 2: **for** $t = 1, \dots, T$ **do**
- 3: compute the sparse codes, *e.g.*, with the homotopy method:
- 4: **for** $i = 1, \dots, n$ **do**
- 5: update α_i by solving (5.8) with $\mathbf{x} = \mathbf{x}_i$ if λ is provided, or (5.10) if μ is provided.
- 6: **end for**
- 7: update the dictionary \mathbf{D} :

$$\mathbf{D} \in \arg \min_{\mathbf{D} \in \mathcal{C}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\alpha_i\|_2^2; \quad (5.23)$$

- 8: **end for**
 - 9: **return** the dictionary \mathbf{D} in \mathcal{C} .
-

appropriate to discard a constraint, solve the resulting unconstrained formulation, before projecting back the solution onto the constraint set. In general, such a procedure can be indeed arbitrarily bad regarding the original constrained problem. In the context of ℓ_0 -regularized dictionary learning, however, such an approach does fortunately make sense. The objective function (5.18) when $\psi = \ell_0$ is indeed invariant to a rescaling operation consisting of multiplying each column of \mathbf{d}_j by a scalar γ_j and the corresponding j -th row of \mathbf{A} by $(1/\gamma_j)$. As a consequence, the minimum value of the objective function is independent of the scale of the columns of \mathbf{D} ; thus the constraint $\mathbf{D} \in \mathcal{C}$ can always be safely ignored and the normalization $\mathbf{D} \leftarrow \Pi_{\mathcal{C}}[\mathbf{D}]$ can be applied at any moment without changing the quality of the dictionary.

In the case of ℓ_1 , the objective function is not invariant to the previous rescaling operation and the MOD update for the dictionary should not be used anymore. Instead, the objective function (5.23) needs to be minimized with another approach. Since it is convex, Lee et al. [2007] propose to use a Newton method in the dual of (5.23), which requires

solving $p \times p$ linear systems at every Newton iteration. Another easy-to-implement choice is the block coordinate descent update of Mairal et al. [2010a], which we present in Algorithm 11. Each step of the block coordinate descent approach is guaranteed to decrease the value of the objective function, while keeping \mathbf{D} in the constraint set \mathcal{C} , and ultimately solve (5.23) since the problem is convex [see Bertsekas, 1999, for convergence results of block coordinate descent methods]. In practice, (5.23) does not need to be solved exactly and performing a few passes, say 10, over the columns of \mathbf{D} seems to perform well in practice.

The update (5.24) can be justified as follows. Minimizing (5.23) with respect to one column \mathbf{d}_j when keeping the other columns fixed can be formulated as

$$\mathbf{d}_j \leftarrow \arg \min_{\mathbf{d} \in \mathbb{R}^m, \|\mathbf{d}\|_2 \leq 1} \left[\sum_{i=1}^n \frac{1}{2} \left\| \mathbf{x}_i - \sum_{l \neq j} \alpha_i[l] \mathbf{d}_l - \alpha_i[j] \mathbf{d} \right\|_2^2 \right].$$

Then, this formulation can be rewritten in a matrix form

$$\mathbf{d}_j \leftarrow \arg \min_{\mathbf{d} \in \mathbb{R}^m, \|\mathbf{d}\|_2 \leq 1} \left[\frac{1}{2} \left\| \mathbf{X} - \mathbf{D}\mathbf{A} + \mathbf{d}_j \boldsymbol{\alpha}^j - \mathbf{d} \boldsymbol{\alpha}^j \right\|_F^2 \right],$$

where \mathbf{d}_j on the right side is the value of the variable \mathbf{d}_j before the update, and $\boldsymbol{\alpha}^j$ in $\mathbb{R}^{1 \times n}$ is the j -th row of the matrix \mathbf{A} . After expanding the Frobenius norm and removing the constant term, we obtain

$$\begin{aligned} \mathbf{d}_j &\leftarrow \arg \min_{\mathbf{d} \in \mathbb{R}^m, \|\mathbf{d}\|_2 \leq 1} \left[-\mathbf{d}^\top (\mathbf{X} - \mathbf{D}\mathbf{A} + \mathbf{d}_j \boldsymbol{\alpha}^j) \boldsymbol{\alpha}^{j\top} + \frac{1}{2} \left\| \mathbf{d} \boldsymbol{\alpha}^j \right\|_F^2 \right] \\ &= \arg \min_{\mathbf{d} \in \mathbb{R}^m, \|\mathbf{d}\|_2 \leq 1} \left[-\mathbf{d}^\top (\mathbf{b}_j - \mathbf{D}\mathbf{c}_j + \mathbf{d}_j \mathbf{C}[j, j]) + \frac{1}{2} \|\mathbf{d}\|_2^2 \mathbf{C}[j, j] \right] \\ &= \arg \min_{\mathbf{d} \in \mathbb{R}^m, \|\mathbf{d}\|_2 \leq 1} \left[\frac{1}{2} \left\| \frac{1}{\mathbf{C}[j, j]} (\mathbf{b}_j - \mathbf{D}\mathbf{c}_j) + \mathbf{d}_j - \mathbf{d} \right\|_2^2 \right], \end{aligned}$$

where the quantities $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_p]$ in $\mathbb{R}^{m \times p}$ and $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_p]$ in $\mathbb{R}^{p \times p}$ are defined in Algorithm 11. Finally, we see from the last equation that updating \mathbf{d}_j amounts to performing one orthogonal projection of the vector $(1/\mathbf{C}[j, j])(\mathbf{b}_j - \mathbf{D}\mathbf{c}_j) + \mathbf{d}_j$ onto the unit Euclidean ball, leading to the update (5.24).

Algorithm 11 Dictionary update with block coordinate descent.

Require: $\mathbf{D}_0 \in \mathcal{C}$ (input dictionary); $\mathbf{X} \in \mathbb{R}^{m \times n}$ (dataset); $\mathbf{A} \in \mathbb{R}^{p \times n}$ (sparse codes);

- 1: Initialization: $\mathbf{D} \leftarrow \mathbf{D}_0$; $\mathbf{B} \leftarrow \mathbf{X}\mathbf{A}^\top$; $\mathbf{C} \leftarrow \mathbf{A}\mathbf{A}^\top$;
- 2: **repeat**
- 3: **for** $j = 1, \dots, p$ **do**
- 4: update the j -th column to optimize for (5.23):

$$\begin{aligned} \mathbf{d}_j &\leftarrow \frac{1}{\mathbf{C}[j, j]}(\mathbf{b}_j - \mathbf{D}\mathbf{c}_j) + \mathbf{d}_j, \\ \mathbf{d}_j &\leftarrow \frac{1}{\max(\|\mathbf{d}_j\|_2, 1)}\mathbf{d}_j. \end{aligned} \tag{5.24}$$

- 5: **end for**
 - 6: **until** convergence;
 - 7: **return** \mathbf{D} (updated dictionary).
-

Block coordinate descent. We have presented in the previous paragraph the alternate minimization method where the dictionary is updated via block coordinate descent. We have also introduced in Section 5.2 a coordinate descent algorithm for dealing with ℓ_1 -penalized problems. It is then natural to combine the two approaches into a block coordinate descent dictionary learning method for (5.17) with $\psi = \ell_1$. We present it in Algorithm 12, where S_λ is the soft-thresholding operator applied element-wise to a vector. The main asset of the block-coordinate descent approach is its simplicity, and ease of implementation. In practice, it can also be very effective.

K-SVD. For ℓ_0 -regularized dictionary learning problems, one of the most popular approach is the K-SVD algorithm of Aharon et al. [2006]. The algorithm is related to the alternate minimization strategy MOD [Engan et al., 1999], but the dictionary update step updates the non-zero coefficients of the matrix \mathbf{A} at the same time. The full approach is detailed in Algorithm 13. The sparse encoding approach is the same as in MOD, but the dictionary elements are updated in a block coordinate fashion, one at a time.

Algorithm 12 Block coordinate descent for (5.17) with $\psi = \ell_1$.

Require: Signals $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$, initial dictionary \mathbf{D}_0 in \mathcal{C} , initial sparse coefficients \mathbf{A}_0 in $\mathbb{R}^{p \times n}$ (possibly equal to zero), regularization parameter λ , number of iterations T .

1: Initialize $\mathbf{D} \leftarrow \mathbf{D}_0$; $\mathbf{A} \leftarrow \mathbf{A}_0$;

2: **for** $t = 1, \dots, T$ **do**

3: update the sparse codes, one row of \mathbf{A} at a time; perform at least one pass (possibly several) of the following iterations:

4: **for** $j = 1, \dots, p$ **do**

5:

$$\boldsymbol{\alpha}^j \leftarrow S_\lambda \left(\boldsymbol{\alpha}^j + \frac{1}{\|\mathbf{d}_j\|_2} \mathbf{d}_j^\top (\mathbf{X} - \mathbf{D}\mathbf{A}) \right); \quad (5.25)$$

6: **end for**

7: update the dictionary \mathbf{D} : perform one pass or more of the inner loop of Algorithm 11.

8: **end for**

9: **return** the dictionary \mathbf{D} in \mathcal{C} .

More precisely, when updating the column \mathbf{d}_j , the algorithm finds the subset of signals Ω that use \mathbf{d}_j in their current sparse decomposition. Then, the non-zero coefficients of the j -th row of \mathbf{A} , denoted by $\boldsymbol{\alpha}_\Omega^j$, are updated at the same time as \mathbf{d}_j in (5.26). Such an update is in fact equivalent to performing a rank-one singular value decomposition, which has inspired the name of the algorithm. A typical technique for solving (5.26) is the power method [see Golub and Van Loan, 2012], which alternates between the optimization of \mathbf{d} with $\boldsymbol{\beta}$ fixed, and of $\boldsymbol{\beta}$ with \mathbf{d} fixed.

Online dictionary learning. Finally, we describe the online dictionary learning method of Mairal et al. [2010a], which we have used in the experiments of this monograph. The approach is based on stochastic approximations, where the goal is to minimize the expectation

$$\min_{\mathbf{D} \in \mathcal{C}} \mathbb{E}_{\mathbf{x}} [L(\mathbf{x}, \mathbf{D})], \quad (5.27)$$

Algorithm 13 K-SVD for (5.18) with $\psi = \ell_0$.

Require: Signals $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ in $\mathbb{R}^{m \times n}$, initial dictionary \mathbf{D}_0 in \mathcal{C} , regularization parameter μ , number of iterations T .

- 1: Initialize $\mathbf{D} \leftarrow \mathbf{D}_0$;
- 2: **for** $t = 1, \dots, T$ **do**
- 3: compute the sparse codes, *e.g.*, with OMP:
- 4: **for** $i = 1, \dots, n$ **do**
- 5:
$$\boldsymbol{\alpha}_i \approx \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}\|_2^2 \quad \text{s.t.} \quad \|\boldsymbol{\alpha}\|_0 \leq \mu;$$
- 6: **end for**
- 7: update the dictionary \mathbf{D} :
- 8: **for** $j = 1, \dots, p$ **do**
- 9: define the set $\Omega \triangleq \{i \in \{1, \dots, n\} : \boldsymbol{\alpha}_i[j] \neq 0\}$;
- 10: update \mathbf{d}_j and the non-zero coefficients of $\boldsymbol{\alpha}^j$:

$$(\mathbf{d}_j, \boldsymbol{\alpha}_\Omega^j) \in \arg \min_{\substack{\mathbf{d} \in \mathbb{R}^m, \|\mathbf{d}\|_2 \leq 1 \\ \beta \in \mathbb{R}^{|\Omega|}}} \left\| \mathbf{X}_\Omega - \sum_{l \neq j} \mathbf{d}_l \boldsymbol{\alpha}_\Omega^l - \mathbf{d} \boldsymbol{\beta}^\top \right\|_F^2; \quad (5.26)$$

- 11: **end for**
 - 12: **end for**
 - 13: **return** the dictionary \mathbf{D} in \mathcal{C} .
-

where L is defined in (5.20), and the expectation is taken with respect to the unknown probability distribution of the data \mathbf{x} . With a discrete probability distribution with a finite sample, (5.27) simply amounts to (5.19).

In a nutshell, the algorithm works by using a stochastic version of the majorization-minimization principle [Mairal, 2013]. It sequentially builds a quadratic surrogate of the expected cost (5.27), which is minimized at every iteration. It can be shown that the quality of approximation of the surrogate near the current estimate improves during the algorithm iterates, and that the approach provides asymptotically a stationary point of (5.27). The basic procedure is summarized in Algorithm 14 and some variants that significantly improves its practical efficiency are presented in [Mairal et al., 2010a].

Assuming that the training set is composed of i.i.d. samples of a data distribution, the inner loop draws one element \mathbf{x}_t at a time, as in the stochastic gradient descent algorithm, and alternates between classical sparse coding steps for computing the decomposition vector $\boldsymbol{\alpha}_t$ of \mathbf{x}_t over the current dictionary \mathbf{D}_{t-1} , and dictionary update steps. The new dictionary \mathbf{D}_t is computed by minimizing over \mathcal{C} the quadratic surrogate (5.28) where the vectors $\boldsymbol{\alpha}_i$ for $i < t$ have been computed during the previous steps of the algorithm. The motivation of the approach is based on two ideas:

- the quadratic surrogate aggregates the past information with a few sufficient statistics, namely the matrices \mathbf{B} and \mathbf{C} . One key component of the convergence analysis shows that the difference between the quadratic surrogate and the objective function at the current estimate converges to zero almost surely.
- with warm restart, the dictionary update is very efficient. In practice, one pass over the dictionary elements seems to perform empirically well.

Variants to improve the efficiency of the online dictionary learning algorithm include: (i) the use of mini-batches instead of drawing a single point at every iteration, (ii) rescaling the past data by updating the matrices $\mathbf{B} \leftarrow \gamma_t \mathbf{B} + \mathbf{x}_t \boldsymbol{\alpha}_t^\top$, with a sequence of parameters γ_t that

converges to one. The implementation of the SPAMS software contains such variants, which have proven to be important in practice [see Mairal et al., 2010a, for additional discussions].

Algorithm 14 Online dictionary learning for (5.17 with $\psi = \ell_1$).

Require: i.i.d signals $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^m$, regularization parameter λ , initial dictionary $\mathbf{D}_0 \in \mathcal{C}$, number of iterations T .

1: Initialization: $\mathbf{B} \in \mathbb{R}^{m \times p} \leftarrow 0$, $\mathbf{C} \in \mathbb{R}^{p \times p} \leftarrow 0$, $\mathbf{D} \leftarrow \mathbf{D}_0$;

2: **for** $t = 1, \dots, T$ **do**

3: draw \mathbf{x}_t and sparsely encode it using homotopy:

$$\boldsymbol{\alpha}_t \in \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{x}_t - \mathbf{D}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1;$$

4: $\mathbf{B} \leftarrow \mathbf{B} + \mathbf{x}_t \boldsymbol{\alpha}_t^\top$;

5: $\mathbf{C} \leftarrow \mathbf{C} + \boldsymbol{\alpha}_t \boldsymbol{\alpha}_t^\top$;

6: update \mathbf{D} using Algorithm 11 (without recomputing the matrices \mathbf{B} and \mathbf{C}) such that

$$\begin{aligned} \mathbf{D} &\leftarrow \arg \min_{\mathbf{D} \in \mathcal{C}} \frac{1}{t} \sum_{i=1}^t \left(\frac{1}{2} \|\mathbf{x}_i - \mathbf{D}\boldsymbol{\alpha}_i\|_2^2 + \lambda \|\boldsymbol{\alpha}_i\|_1 \right), \\ &= \arg \min_{\mathbf{D} \in \mathcal{C}} \frac{1}{t} \left(\frac{1}{2} \text{trace}(\mathbf{D}^\top \mathbf{D} \mathbf{C}) - \text{trace}(\mathbf{D}^\top \mathbf{B}) \right); \end{aligned} \quad (5.28)$$

7: **end for**

8: **return** learned dictionary \mathbf{D} in \mathcal{C} .

Note that along the same lines, a similar algorithm called “recursive least square dictionary learning” was independently developed by Skretting and Engan [2010] focusing on the ℓ_0 -penalty.

Other extensions. Because the literature on dictionary learning algorithms is vast, we have presented a few ones that address the formulations (5.17) or (5.18). Extending these approaches to other settings, such as other loss functions than the square loss, or other regularization functions, is easy in general by replacing the sparse encoding and dictionary update steps by ad hoc optimization techniques.

5.6 Other optimization techniques

For simplicity, we have focused on sparse estimation problems involving least-square cost functions, and have left aside structured sparsity penalties such as the group Lasso. Some of the methods we have described can be easily extended to these settings, in particular the proximal gradient and coordinate descent algorithms are relatively flexible. We leave these extensions to more exhaustive reviews [see Bach et al., 2012a, for instance].

Other classes of algorithms are also popular for solving convex sparse estimation problems. For example, proximal splitting techniques [Combettes and Pesquet, 2011, Chambolle and Pock, 2011, Condat, 2013, Raguet et al., 2013], such as the alternating direction method of multipliers [see Boyd et al., 2011], are often used in signal processing. These methods are very flexible, simple to implement, and can be applied to a large range of problems. However, they require tuning (at least) one parameter to obtain fast convergence in practice, which, to the best of our knowledge, cannot be done automatically. In the context of machine learning, stochastic optimization techniques have also been developed [*e.g.* Duchi et al., 2011, Xiao, 2010]. These methods are adapted to large-scale problems; even though, they may seem simple at first sight, they require a careful choice of a step size parameter and thus remain a bit difficult to use.

Finally, we have also omitted quasi-Newton methods such as variants of L-BFGS [Schmidt et al., 2011], which can exploit the curvature of the objective function to achieve fast convergence.

6

Conclusions

In this monograph, we have reviewed a large number of applications of dictionary learning in image processing and computer vision and presented basic sparse estimation tools. We started with a historical tour of sparse estimation in signal processing and statistics, before moving to more recent concepts such as sparse recovery and dictionary learning. Subsequently, we have shown that dictionary learning is related to matrix factorization techniques, and that it is particularly effective for modeling natural image patches. As a consequence, it has been used for tackling several image processing problems and is a key component of many state-of-the-art methods in visual recognition. We have concluded the monograph with a presentation of optimization techniques that should make dictionary learning easy to use for researchers that are not experts of the field.

With efficient software available, it is remarkable to see that the interest for dictionary learning and sparse estimation is still acute today, and seems to increase in other scientific communities than image processing or computer vision. In particular, new applications of sparse matrix factorization are found in neuroscience [Varoquaux et al., 2011], bioinformatics for processing gene expression data [see Barillot et al.,

2012, Section 5.4], audio processing [Baby et al., 2014], and in astrophysics [Vinci et al., 2014]. We have also put the emphasis on natural images, and we have omitted the application of dictionary learning to other domains, such as hyperspectral data, where significant success have been obtained [Charles et al., 2011, Xing et al., 2012], or to disparity maps for stereo imaging [Tosic et al., 2011].

Exploring the possible use of dictionary learning in these fields would be of great interest and high importance, but unfortunately this goal falls beyond the scope of this monograph.

Acknowledgements

This monograph was partially supported by the European Research Council (SIERRA and VideoWorld projects), the project Gargantua funded by the program Mastodons-CNRS, the Microsoft Research-Inria joint centre, and a French grant from the Agence Nationale de la Recherche (MACARON project, ANR-14-CE23-0003-01).

The authors would like to thank Ori Bryt and Michael Elad for providing the output of their face compression algorithm (see Figure 3.9). Julien Mairal would also like to thank Zaid Harchaoui, Florent Perronnin and Adrien Gaidon for useful discussions leading to improvements of this monograph.

References

- J. Abernethy, F. Bach, T. Evgeniou, and J.-P. Vert. A new approach to collaborative filtering: operator estimation with spectral regularization. *Journal of Machine Learning Research*, 10:803–826, 2009.
- A. Agarwal, A. Anandkumar, P. Jain, P. Netrapalli, and R. Tandon. Learning sparsely used overcomplete dictionaries via alternating minimization. *preprint arXiv:1310.7991*, 2013.
- M. Aharon, M. Elad, and A. Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322, 2006.
- N. Ahmed, T. Natarajan, and K. R. Rao. Discrete cosine transform. *IEEE Transactions on Computers*, 100(1):90–93, 1974.
- H. Akaike. Information theory and an extension of the maximum likelihood principle. In *Second International Symposium on Information Theory*, volume 1, pages 267–281, 1973.
- S. Arora, R. Ge, and A. Moitra. New algorithms for learning incoherent and overcomplete dictionaries. In *Proceedings of the Annual Conference on Computational Learning Theory*, 2014.
- S. P. Awate and R. T. Whitaker. Unsupervised, information-theoretic, adaptive image filtering for image restoration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(3):364–376, 2006.

- D. Baby, T. Virtanen, T. Barker, and H. Van hamme. Coupled dictionary training for exemplar-based speech enhancement. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2014.
- F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Optimization with sparsity-inducing penalties. *Foundation and Trends in Machine Learning*, 4:1–106, 2012a.
- F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. Structured sparsity through convex optimization. *Statistical Science*, 27(4):450–468, 2012b.
- S. Bakin. *Adaptive regression and model selection in data mining problems*. PhD thesis, 1999.
- R. G. Baraniuk, V. Cevher, M. Duarte, and C. Hegde. Model-based compressive sensing. *IEEE Transactions on Information Theory*, 56(4):1982–2001, 2010.
- E. Barillot, L. Calzone, P. Hupe, J.-P. Vert, and A. Zinovyev. *Computational systems biology of cancer*. CRC Press, 2012.
- A. Barron, J. Rissanen, and B. Yu. The minimum description length principle in coding and modeling. *IEEE Transactions on Information Theory*, 44(6):2743–2760, 1998.
- M. J. Bayarri and J. O Berger. The interplay of bayesian and frequentist analysis. *Statistical Science*, 19(1):58–80, 2004.
- 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.
- M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.
- M. Belkin and P. Niyogi. Semi-supervised learning on Riemannian manifolds. *Machine learning*, 56(1-3):209–239, 2004.
- A. J. Bell and T. J. Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural computation*, 7(6):1129–1159, 1995.
- A. J. Bell and T. J. Sejnowski. The “independent components” of natural scenes are edge filters. *Vision Research*, 37(23):3327–3338, 1997.
- Y. Bengio. Learning deep architectures for AI. *Foundations and Trends in Machine Learning*, 2(1):1–127, 2009.
- M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester. Image inpainting. In *Proceedings of the ACM SIGGRAPH Conference*, 2000.

- D. P. Bertsekas. *Nonlinear programming*. Athena Scientific, 1999. 2nd edition.
- P. J. Bickel. Parametric robustness: small biases can be worthwhile. *The Annals of Statistics*, 12(4):864–879, 1984.
- C. M. Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- T. Blumensath and M. E. Davies. Iterative hard thresholding for compressed sensing. *Applied and Computational Harmonic Analysis*, 27(3):265–274, 2009.
- L. Bo, X. Ren, and D. Fox. Hierarchical matching pursuit for image classification: Architecture and fast algorithms. In *Advances in Neural Information Processing Systems (NIPS)*, 2011.
- L. Bo, X. Ren, and D. Fox. Multipath sparse coding using hierarchical matching pursuit. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- J.M. Borwein and A.S. Lewis. *Convex analysis and nonlinear optimization: theory and examples*. Springer, 2006.
- T. Bossomaier and A. W. Snyder. Why spatial frequency processing in the visual cortex? *Vision Research*, 26(8):1307–1309, 1986.
- L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2008.
- Y-L. Boureau, F. Bach, Y. LeCun, and J. Ponce. Learning mid-level features for recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- Y-L. Boureau, N. Le Roux, F. Bach, J. Ponce, and Y. LeCun. Ask the locals: Multi-way local pooling for image recognition. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.
- D. M. Bradley and J. A. Bagnell. Differential sparse coding. In *Advances in Neural Information Processing Systems (NIPS)*, 2008.
- P. Brucker. An $O(n)$ algorithm for quadratic knapsack problems. *Operations Research Letters*, 3(3):163–166, 1984.

- J. Bruna and S. Mallat. Invariant scattering convolution networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1872–1886, 2013.
- O. Bryt and M. Elad. Compression of facial images using the K-SVD algorithm. *Journal of Visual Communication and Image Representation*, 19(4):270–282, 2008.
- A. Buades, B. Coll, and J.-M. Morel. A review of image denoising algorithms, with a new one. *SIAM Journal on Multiscale Modeling and Simulation*, 4(2):490–530, 2005.
- A. Buades, B. Coll, J.-M. Morel, and C. Sbert. Self-similarity driven color demosaicking. *IEEE Transactions on Image Processing*, 18(6):1192–1202, 2009.
- P. Bühlmann and S. Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer, 2011.
- F. Bunea, A. B. Tsybakov, and M. H. Wegkamp. Aggregation for Gaussian regression. *The Annals of Statistics*, 35(4):1674–1697, 2007.
- T. T. Cai. Adaptive wavelet estimation: a block thresholding and oracle inequality approach. *Annals of Statistics*, 27(3):898–924, 1999.
- E. J. Candès and D. L. Donoho. Recovering edges in ill-posed inverse problems: Optimality of curvelet frames. *Annals of Statistics*, 30(3):784–842, 2002.
- E. J. Candès and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005.
- E. J. Candès and M. B. Wakin. An introduction to compressive sampling. *IEEE Signal Processing Magazine*, 25(2):21–30, 2008.
- E. J. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006.
- E.J. Candès, M. Wakin, and S. Boyd. Enhancing sparsity by reweighted ℓ_1 minimization. *Journal of Fourier Analysis and Applications*, 14(5):877–905, 2008.
- M. Carandini, J. B. Demb, V. Mante, D. J. Tolhurst, Y. Dan, B. A. Olshausen, J. L. Gallant, and N. C. Rust. Do we know what the early visual system does? *The Journal of Neuroscience*, 25(46):10577–10597, 2005.
- J.-F. Cardoso. Dependence, correlation and Gaussianity in independent component analysis. *Journal of Machine Learning Research*, 4:1177–1203, 2003.

- A. Castrodad and G. Sapiro. Sparse modeling of human actions from motion imagery. *International Journal of Computer Vision*, 100(1):1–15, 2012.
- A. Chambolle. Total variation minimization and a class of binary MRF models. In *Proceedings of the 5th International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition*, 2005.
- A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, 2011.
- S. G. Chang, B. Yu, and M. Vetterli. Spatially adaptive wavelet thresholding with context modeling for image denoising. *IEEE Transactions on Image Processing*, 9(9):1522–1531, 2000a.
- S. G. Chang, B. Yu, and M. Vetterli. Adaptive wavelet thresholding for image denoising and compression. *IEEE Transactions on Image Processing*, 9(9):1532–1546, 2000b.
- A. S. Charles, B. A. Olshausen, and C. J. Rozell. Learning sparse codes for hyperspectral imagery. *IEEE Journal of Selected Topics in Signal Processing*, 5(5):963–978, 2011.
- P. Chatterjee and P. Milanfar. Patch-based near-optimal image denoising. *IEEE Transactions on Image Processing*, 21(4):1635–1649, 2012.
- S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Scientific Computing*, 20:33–61, 1999.
- Y. Chen, J. Mairal, and Z. Harchaoui. Fast and robust archetypal analysis for representation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- E. C. Chi, H. Zhou, G. K. Chen, D. O. Del Vecchio, and K. Lange. Genotype imputation via matrix completion. *Genome research*, 23(3):509–518, 2013.
- J. F. Claerbout and F. Muir. Robust modeling with erratic data. *Geophysics*, 38(5):826–844, 1973.
- A. Coates, A. Y. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *International Conference on Artificial Intelligence and Statistics*, 2011.
- R. R. Coifman and D. L. Donoho. Translation-invariant de-noising. In *Lecture Notes in Statistics*, volume 103, pages 125–150. 1995.
- P. L. Combettes and J.-C. Pesquet. *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, chapter Proximal Splitting Methods in Signal Processing. Springer, 2011.

- L. Condat. A primal–dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms. *Journal of Optimization Theory and Applications*, 158(2):460–479, 2013.
- S. F. Cotter, J. Adler, B. Rao, and K. Kreutz-Delgado. Forward sequential algorithms for best basis selection. In *IEEE Proceedings of Vision Image and Signal Processing*, pages 235–244, 1999.
- F. Couzinie-Devy, J. Mairal, F. Bach, and J. Ponce. Dictionary learning for deblurring and digital zoom. *preprint arXiv:1110.0957*, 2011.
- T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley and Sons, 2006. 2nd edition.
- A. Criminisi, P. Pérez, and K. Toyama. Region filling and object removal by exemplar-based image inpainting. *IEEE Transactions on Image Processing*, 13(9):1200–1212, 2004.
- M. S. Crouse, R. D. Nowak, and R. G. Baraniuk. Wavelet-based statistical signal processing using hidden Markov models. *IEEE Transactions on Signal Processing*, 46(4):886–902, 1998.
- G. Csurka, C. Dance, L. Fan, J. Willamowski, and C. Bray. Visual categorization with bags of keypoints. In *Proceedings of the workshop on statistical learning in computer vision, ECCV*, 2004.
- A. Cutler and L. Breiman. Archetypal analysis. *Technometrics*, 36(4):338–347, 1994.
- K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. Image denoising by sparse 3D transform-domain collaborative filtering. *IEEE Transactions on Image Processing*, 16(8):2080–2095, 2007a.
- K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. Color image denoising via sparse 3D collaborative filtering with grouping constraint in luminance-chrominance space. In *IEEE International Conference on Image Processing (ICIP)*, 2007b.
- K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. BM3D image denoising with shape-adaptive principal component analysis. In *SPARS’09-Signal Processing with Adaptive Sparse Structured Representations*, 2009.
- N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005.
- A. d’Aspremont and L. El Ghaoui. Testing the nullspace property using semidefinite programming. *Mathematical Programming*, 127(1):123–144, 2011.

- I. Daubechies. Orthonormal bases of compactly supported wavelets. *Communications on pure and applied mathematics*, 41(7):909–996, 1988.
- I. Daubechies, M. Defrise, and C. De Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57(11):1413–1457, 2004.
- I. Daubechies, R. DeVore, M. Fornasier, and C. S. Güntürk. Iteratively reweighted least squares minimization for sparse recovery. *Communications on Pure and Applied Mathematics*, 63(1):1–38, 2010.
- J. G. Daugman. Uncertainty relation for resolution in space, spatial frequency, and orientation optimized by two-dimensional visual cortical filters. *Journal of the Optical Society of America A*, 2(7):1160–1169, 1985.
- A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society Series B*, 39(1):1–38, 1977.
- M. N. Do and M. Vetterli. The contourlet transform: an efficient directional multiresolution image representation. *IEEE Transactions on Image Processing*, 14(12):2091–2106, 2005.
- W. Dong, L. Zhang, and G. Shi. Centralized sparse representation for image restoration. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011a.
- W. Dong, L. Zhang, G. Shi, and X. Wu. Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization. *IEEE Transactions on Image Processing*, 20(7):1838–1857, 2011b.
- W. Dong, L. Zhang, G. Shi, and X. Li. Nonlocally centralized sparse representation for image restoration. *IEEE Transactions on Image Processing*, 22(4):1620–1630, 2013.
- D. L. Donoho. Wedgelets: Nearly minimax estimation of edges. *Annals of Statistics*, 27(3):859–897, 1999.
- D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006.
- D. L. Donoho and I. M. Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *Journal of the American Statistical Association*, 90(432):1200–1224, 1995.
- D. L. Donoho and J. M. Johnstone. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81(3):425–455, 1994.

- J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the l_1 -ball for learning in high dimensions. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.
- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for on-line learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- C. Eckart and G. Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- B. Efron and C. Morris. Limiting the risk of Bayes and empirical Bayes estimators — part I: the Bayes case. *Journal of the American Statistical Association*, 66(336):807–815, 1971.
- B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32(2):407–499, 2004.
- A. A. Efros and T. K. Leung. Texture synthesis by non-parametric sampling. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 1999.
- M. A. Efronson. Multiple regression analysis. *Mathematical methods for digital computers*, 9(1):191–203, 1960.
- M. B. Eisen, P. T. Spellman, P. O. Brown, and D. Botstein. Cluster analysis and display of genome-wide expression patterns. *Proceedings of the National Academy of Sciences of the United States of America*, 95(25):14863–14868, 1998.
- M. Elad. *Sparse and redundant representations: from theory to applications in signal and image processing*. Springer, 2010.
- 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.
- M. Elad and A. M. Bruckstein. A generalized uncertainty principle and sparse representation in pairs of bases. *IEEE Transactions on Information Theory*, 48(9):2558–2567, 2002.
- E. Elhamifar and R. Vidal. Sparse subspace clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- E. Elhamifar, G. Sapiro, and R. Vidal. See all by looking at a few: Sparse modeling for finding representative objects. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

- K. Engan, S. O. Aase, and J. H. Husoy. Method of optimal directions for frame design. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 1999.
- J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- O. D. Faugeras. Digital color image processing within the framework of a human visual model. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 27(4):380–393, 1979.
- M. Fazel. *Matrix rank minimization with applications*. PhD thesis, Stanford University, 2002.
- M. Fazel, H. Hindi, and S. Boyd. A rank minimization heuristic with application to minimum order system approximation. In *American Control Conference*, volume 6, pages 4734–4739, 2001.
- M. Fazel, H. Hindi, and S. P. Boyd. Log-det heuristic for matrix rank minimization with applications to Hankel and Euclidean distance matrices. In *Proceedings of the American Control Conference*, volume 3, pages 2156–2162, 2003.
- 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):793–830, 2009.
- D. J. Field. Relations between the statistics of natural images and the response properties of cortical cells. *Journal of the Optical Society of America A*, 4(12):2379–2394, 1987.
- M. A. T. Figueiredo and R. D. Nowak. An EM algorithm for wavelet-based image restoration. *IEEE Transactions on Image Processing*, 12(8):906–916, 2003.
- M. A. T. Figueiredo and R. D. Nowak. A bound optimization approach to wavelet-based image deconvolution. In *Proceedings of the International Conference on Image Processing (ICIP)*, 2005.
- M. A. T. Figueiredo, J. M. Bioucas-Dias, and R. D. Nowak. Majorization-minimization algorithms for wavelet-based image restoration. *IEEE Transactions on Image Processing*, 16(12):2980–2991, 2007.
- R. W. Floyd and L. Steinberg. An adaptive algorithm for spatial grey scale. In *Proceedings of the Society of Information Display*, volume 17, pages 75–77, 1976.
- D. A. Forsyth and J. Ponce. *Computer vision: a modern approach*. Prentice Hall, 2012. 2nd edition.

- I. E. Frank and J. H. Friedman. A statistical view of some chemometrics regression tools. *Technometrics*, 35(2):109–135, 1993.
- W. T. Freeman and E. H. Adelson. The design and use of steerable filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(9):891–906, 1991.
- W. T. Freeman, T. R. Jones, and E. C. Pasztor. Example-based super-resolution. *IEEE Computer Graphics and Applications*, 22(2):56–65, 2002.
- J. H. Friedman and W. Stuetzle. Projection pursuit regression. *Journal of the American Statistical Association*, 76(376):817–823, 1981.
- W.J. Fu. Penalized regressions: the bridge versus the lasso. *Journal of Computational and Graphical Statistics*, 7(3):397–416, 1998.
- J. J. Fuchs. Recovery of exact sparse representations in the presence of bounded noise. *IEEE Transactions on Image Processing*, 51(10):3601–3608, 2005.
- G. M. Furnival and R. W. Wilson. Regressions by leaps and bounds. *Technometrics*, 16(4):499–511, 1974.
- D. Gabor. Theory of communication. *Journal of the Institution of Electrical Engineers*, 93(26):429–441, 1946.
- S. Gao, I. W.-H. Tsang, L.-T. Chia, and P. Zhao. Local features are not lonely—Laplacian sparse coding for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- S. Gao, W.-H. Tsang, and L.-T. Chia. Laplacian sparse coding, hypergraph Laplacian sparse coding, and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1):92–104, 2013.
- P. Garrigues and B. A. Olshausen. Group sparse coding with a Laplacian scale mixture prior. In *Advances in Neural Information Processing Systems (NIPS)*, 2010.
- G. Gasso, A. Rakotomamonjy, and S. Canu. Recovering sparse signals with non-convex penalties and DC programming. *IEEE Transactions on Signal Processing*, 57(12):4686–4698, 2009.
- Q. Geng, H. Wang, and J. Wright. On the local correctness of ℓ_1 -minimization for dictionary learning. *preprint arXiv:1101.5672*, 2011.
- A. Gersho and R. M. Gray. *Vector quantization and signal compression*. Kluwer Academic Publishers, 1992.
- M. Gharavi-Alkhansari and T. S. Huang. A fast orthogonal matching pursuit algorithm. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 1998.

- R. Giryes and M. Elad. Sparsity based Poisson denoising with dictionary learning. *IEEE Transactions on Image Processing*, 2014. to appear.
- D. Glasner, S. Bagon, and M. Irani. Super-resolution from a single image. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2009.
- G. H. Golub and C. F. Van Loan. *Matrix computations*, volume 3. Johns Hopkins University Press, 2012. 4th edition.
- Y. Grandvalet and S. Canu. Outcomes of the equivalence of adaptive ridge with least absolute shrinkage. In *Advances in Neural Information Processing Systems (NIPS)*, 1999.
- H. Grassmann. LXXXVII. On the theory of compound colours. *Philosophical Magazine Series 4*, 7(45):254–264, 1854.
- K. Grauman and T. Darrell. The pyramid match kernel: Discriminative classification with sets of image features. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2005.
- E. Greenshtein. Best subset selection, persistence in high-dimensional statistical learning and optimization under ℓ_1 constraint. *The Annals of Statistics*, 34(5):2367–2386, 2006.
- K. Gregor and Y. LeCun. Learning fast approximations of sparse coding. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2010.
- K. Gregor, A. Szlam, and Y. LeCun. Structured sparse coding via lateral inhibition. In *Advances in Neural Information Processing Systems (NIPS)*, 2011.
- R. Gribonval and M. Nielsen. Sparse representations in unions of bases. *IEEE Transactions on Information Theory*, 49(12):3320–3325, 2003.
- R. Gribonval and K. Schnass. Dictionary identification–sparse matrix-factorization via ℓ_1 -minimization. *IEEE Transactions on Information Theory*, 56(7):3523–3539, 2010.
- R. Gribonval, V. Cevher, and M. E. Davies. Compressible distributions for high-dimensional statistics. *IEEE Transactions on Information Theory*, 58(8):5016–5034, 2012.
- R. Gribonval, R. Jenatton, F. Bach, M. Kleinsteuber, and M. Seibert. Sample complexity of dictionary learning and other matrix factorizations. *preprint arXiv:1312.3790*, 2013.
- R. Gribonval, R. Jenatton, and F. Bach. Sparse and spurious: dictionary learning with noise and outliers. *preprint arXiv:1407.5155*, 2014.

- T. Guha and R. K. Ward. Learning sparse representations for human action recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(8):1576–1588, 2012.
- A. Haar. Zur theorie der orthogonalen funktionensysteme. *Mathematische Annalen*, 69(3):331–371, 1910.
- P. Hall, G. Kerkycharian, and D. Picard. On the minimax optimality of block thresholded wavelet estimators. *Statistica Sinica*, 9(1):33–49, 1999.
- T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning*. Springer, 2009. 2nd edition.
- S. Hawe, M. Seibert, and M. Kleinsteuber. Separable dictionary learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- J. Héroult, C. Jutten, and B. Ans. Détection de grandeurs primitives dans un message composite par une architecture de calcul neuromimétique en apprentissage non supervisé. In *Actes du Xème Colloque GRETSI*, 1985.
- K. K. Herrity, A. C. Gilbert, and J. A. Tropp. Sparse approximation via iterative thresholding. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing, (ICASSP)*, 2006.
- G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800, 2002.
- R. R. Hocking. A Biometrics invited paper. The analysis and selection of variables in linear regression. *Biometrics*, 32:1–49, 1976.
- A. E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- H. Hotelling. Analysis of a complex of statistical variables into principal components. *Journal of educational psychology*, 24(6):417, 1933.
- J. Huang and T. Zhang. The benefit of group sparsity. *Annals of Statistics*, 38(4):1978–2004, 2010.
- J. Huang, Z. Zhang, and D. Metaxas. Learning with structured sparsity. *Journal of Machine Learning Research*, 12:3371–3412, 2011.
- K. Huang and S. Aiyente. Sparse representation for signal classification. In *Advances in Neural Information Processing Systems (NIPS)*, 2006.
- D. H. Hubel and T. N. Wiesel. Receptive fields and functional architecture of monkey striate cortex. *The Journal of physiology*, 195(1):215–243, 1968.
- A. Hyvärinen, P. O. Hoyer, and M. Inki. Topographic independent component analysis. *Neural computation*, 13(7):1527–1558, 2001.

- A. Hyvärinen, J. Karhunen, and E. Oja. *Independent component analysis*. John Wiley and Sons, 2004.
- A. Hyvärinen, J. Hurri, and P. O. Hoyer. *Natural Image Statistics: A Probabilistic Approach to Early Computational Vision*. Springer, 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.
- K. Jarrett, K. Kavukcuoglu, M. Ranzato, and Y. LeCun. What is the best multi-stage architecture for object recognition? In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2009.
- H. Jégou, M. Douze, and C. Schmid. Improving bag-of-features for large scale image search. *International Journal of Computer Vision*, 87(3):316–336, 2010.
- R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for sparse hierarchical dictionary learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2010a.
- R. Jenatton, G. Obozinski, and F. Bach. Structured sparse principal component analysis. In *Proceedings of International Workshop on Artificial Intelligence and Statistics (AISTATS)*, 2010b.
- R. Jenatton, J.-Y. Audibert, and F. Bach. Structured variable selection with sparsity-inducing norms. *Journal of Machine Learning Research*, 12:2777–2824, 2011a.
- R. Jenatton, J. Mairal, G. Obozinski, and F. Bach. Proximal methods for hierarchical sparse coding. *Journal of Machine Learning Research*, 12:2297–2334, 2011b.
- A. Juditsky and A. Nemirovski. On verifiable sufficient conditions for sparse signal recovery via ℓ_1 -minimization. *Mathematical Programming*, 127(1): 57–88, 2011.
- V. Katkovnik, A. Foi, K. Egiazarian, and J. Astola. From local kernel to non-local multiple-model image denoising. *International Journal of Computer Vision*, 86(1):1–32, 2010.
- K. Kavukcuoglu, M. Ranzato, R. Fergus, and Y. LeCun. Learning invariant features through topographic filter maps. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- K. Kavukcuoglu, M. Ranzato, and Y. LeCun. Fast inference in sparse coding algorithms with applications to object recognition. *preprint arXiv:1010.3467*, 2010a.

- K. Kavukcuoglu, P. Sermanet, Y-L. Boureau, K. Gregor, M. Mathieu, and Y. LeCun. Learning convolutional feature hierarchies for visual recognition. In *Advances in Neural Information Processing Systems (NIPS)*, 2010b.
- K. N. Kay, T. Naselaris, R. J. Prenger, and J. L. Gallant. Identifying natural images from human brain activity. *Nature*, 452(7185):352–355, 2008.
- C. Kervrann and J. Boulanger. Optimal spatial adaptation for patch-based image denoising. *IEEE Transactions on Image Processing*, 15(10):2866–2878, 2006.
- R. Kimmel. Demosaicing: image reconstruction from color ccd samples. *IEEE Transactions on Image Processing*, 8(9):1221–1228, 1999.
- J. Koenderink and A. Van Doorn. The structure of locally orderless images. *International Journal of Computer Vision*, 31(2/3):159–168, 1999.
- P. Koniusz and K. Mikolajczyk. Spatial coordinate coding to reduce histogram representations, dominant angle and colour pyramid match. In *Proceedings of the International Conference on Image Processing (ICIP)*, 2011.
- P. Koniusz, F. Yan, and K. Mikolajczyk. Comparison of mid-level feature coding approaches and pooling strategies in visual concept detection. *Computer Vision and Image Understanding*, 117(5):479–492, 2013.
- H. J. Kushner and G. Yin. *Stochastic approximation and recursive algorithms and applications*. Springer, 2003.
- K. Lange, D. R. Hunter, and I. Yang. Optimization transfer using surrogate objective functions. *Journal of Computational and Graphical Statistics*, 9(1):1–20, 2000.
- I. Laptev. On space-time interest points. *International Journal of Computer Vision*, 64(2-3):107–123, 2005.
- 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.
- E. Le Pennec and S. Mallat. Sparse geometric image representations with bandelets. *IEEE Transactions on Image Processing*, 14(4):423–438, 2005.
- Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998a.
- Y. LeCun, L. Bottou, G. Orr, and K. Muller. Efficient backprop. In G. Orr and Muller K., editors, *Neural Networks: Tricks of the trade*. Springer, 1998b.

- D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–791, 1999.
- H. Lee, A. Battle, R. Raina, and A.Y. Ng. Efficient sparse coding algorithms. In *Advances in Neural Information Processing Systems (NIPS)*, 2007.
- A. Levin, B. Nadler, F. Durand, and W. T. Freeman. Patch complexity, finite pixel correlations and optimal denoising. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2012.
- M. S. Lewicki. Efficient coding of natural sounds. *Nature neuroscience*, 5(4):356–363, 2002.
- M. S. Lewicki and B. A. Olshausen. Probabilistic framework for the adaptation and comparison of image codes. *Journal of the Optical Society of America A*, 16(7):1587–1601, 1999.
- M. S. Lewicki and T. J. Sejnowski. Learning overcomplete representations. *Neural computation*, 12(2):337–365, 2000.
- C.-K. Li and W. So. Isometries of the ℓ_p norm. *The American Mathematical Monthly*, 101(5):452–453, 1994.
- Y. Li and D. P. Huttenlocher. Sparse long-range random field and its application to image denoising. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2008.
- M. S. Livingstone and D. H. Hubel. Anatomy and physiology of a color system in the primate visual cortex. *The Journal of Neuroscience*, 4(1):309–356, 1984.
- D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2):91–110, 2004.
- B. D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of the 7th Joint Conference on Artificial Intelligence (IJCAI)*, 1981.
- D. J. C. MacKay. *Information theory, inference, and learning algorithms*. Cambridge university press, 2003.
- M. Mahmoudi and G. Sapiro. Fast image and video denoising via nonlocal means of similar neighborhoods. *IEEE Signal Processing Letters*, 12(12):839–842, 2005.
- J. Mairal. *Sparse coding for machine learning, image processing and computer vision*. PhD thesis, Ecole Normale Supérieure de Cachan, 2010. <http://tel.archives-ouvertes.fr/tel-00595312>.

- J. Mairal. Stochastic majorization-minimization algorithms for large-scale optimization. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- J. Mairal and B. Yu. Complexity analysis of the Lasso regularization path. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2012.
- J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Discriminative learned dictionaries for local image analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008a.
- J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Supervised dictionary learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2008b.
- J. Mairal, M. Elad, and G. Sapiro. Sparse representation for color image restoration. *IEEE Transactions on Image Processing*, 17(1):53–69, 2008c.
- J. Mairal, M. Leordeanu, F. Bach, M. Hebert, and J. Ponce. Discriminative sparse image models for class-specific edge detection and image interpretation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2008d.
- J. Mairal, G. Sapiro, and M. Elad. Learning multiscale sparse representations for image and video restoration. *SIAM Multiscale Modeling and Simulation*, 7(1):214–241, 2008e.
- J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Non-local sparse models for image restoration. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2009.
- J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11:19–60, 2010a.
- J. Mairal, R. Jenatton, G. Obozinski, and F. Bach. Network flow algorithms for structured sparsity. In *Advances in Neural Information Processing Systems (NIPS)*, 2010b.
- 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.
- J. Mairal, F. Bach, and J. Ponce. Task-driven dictionary learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4):791–804, 2012.
- S. Mallat. *A wavelet tour of signal processing*. Academic press, 2008. 3rd edition.

- S. Mallat and Z. Zhang. Matching pursuit in a time-frequency dictionary. *IEEE Transactions on Signal Processing*, 41(12):3397–3415, 1993.
- S. G. Mallat. A theory for multiresolution signal decomposition: the wavelet representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(7):674–693, 1989.
- C. L. Mallows. Choosing variables in a linear regression: A graphical aid. unpublished paper presented at the Central Regional Meeting of the Institute of Mathematical Statistics, Manhattan, Kansas, 1964.
- C. L. Mallows. Choosing a subset regression. unpublished paper presented at the Joint Statistical Meeting, Los Angeles, California, 1966.
- H. Markowitz. Portfolio selection. *Journal of Finance*, 7(1):77–91, 1952.
- D. Martin, C. Fowlkes, Doron Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2001.
- D. R. Martin, C. C. Fowlkes, and J. Malik. Learning to detect natural image boundaries using local brightness, color, and texture cues. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(5):530–549, 2004.
- P. Massart. *Concentration Inequalities and Model Selection: Ecole d’été de Probabilités de Saint-Flour 23*. Springer, 2003.
- I. Matthews, T. Ishikawa, and S. Baker. The template update problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(6):810–815, 2004.
- A. Maurer and M. Pontil. k -dimensional coding schemes in Hilbert spaces. *IEEE Transactions on Information Theory*, 56(11):5839–5846, 2010.
- J. C. Maxwell. On the theory of compound colours, and the relations of the colours of the spectrum. *Philosophical Transactions of the Royal Society of London*, pages 57–84, 1860.
- X. Mei and H. Ling. Robust visual tracking using ℓ_1 -minimization. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2009.
- N. Meinshausen and P. Bühlmann. Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4):417–473, 2010.
- D. Menon and G. Calvagno. Color image demosaicking: an overview. *Signal Processing: Image Communication*, 26(8):518–533, 2011.
- C. A. Micchelli and M. Pontil. Learning the kernel function via regularization. In *Journal of Machine Learning Research*, volume 6, pages 1099–1125, 2005.

- K. Mikolajczyk and C. Schmid. A performance evaluation of local descriptors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(10):1615–1630, 2005.
- N. Murray and F. Perronnin. Generalized max pooling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- E. A. Nadaraya. On estimating regression. *Theory of Probability and its Applications*, 9:141, 1964.
- I. Naseem, R. Togneri, and M. Bennamoun. Linear regression for face recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(11):2106–2112, 2010.
- N. M. Nasrabadi and R. A. King. Image coding using vector quantization: A review. *IEEE Transactions on Communications*, 36(8):957–971, 1988.
- B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM Journal on Computing*, 24:227–234, 1995.
- J. Nathans, D. Thomas, and D. S. Hogness. Molecular genetics of human color vision: the genes encoding blue, green, and red pigments. *Science*, 232(4747):193–202, 1986.
- Y. Nesterov. Gradient methods for minimizing composite objective function. *Mathematical Programming*, 140(1):125–161, 2013.
- I. Newton. Hypothesis explaining the properties of light. In *The History of the Royal Society*, volume 3, pages 247–269. T. Birch, 1675. text published in 1757.
- S. Nishimoto, A. T. Vu, T. Naselaris, Y. Benjamini, B. Yu, and J. L. Gallant. Reconstructing visual experiences from brain activity evoked by natural movies. *Current Biology*, 21(19):1641–1646, 2011.
- R. D. Nowak and M. A. T. Figueiredo. Fast wavelet-based image deconvolution using the EM algorithm. In *Conference Record of the Thirty-Fifth Asilomar Conference on Signals, Systems and Computers.*, 2001.
- G. Obozinski, B. Taskar, and M. I. Jordan. Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, 20(2):231–252, 2009.
- T. Ojala, M. Pietikainen, and T. Maenpaa. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(7):971–987, 2002.

- B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381:607–609, 1996.
- B. A. Olshausen and D. J. Field. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, 37(23):3311–3325, 1997.
- B. A. Olshausen and D. J. Field. How close are we to understanding V1? *Neural computation*, 17(8):1665–1699, 2005.
- M. R. Osborne, B. Presnell, and B. A. Turlach. A new approach to variable selection in least squares problems. *IMA journal of numerical analysis*, 20(3):389–403, 2000.
- P. Paatero and U. Tapper. Positive matrix factorization: a non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126, 1994.
- T. Park and G. Casella. The Bayesian Lasso. *Journal of the American Statistical Association*, 103(482):681–686, 2008.
- Y. C. Pati, R. Rezaiifar, and P. S. Krishnaprasad. Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Proceedings of The Twenty-Seventh Asilomar Conference on Signals, Systems and Computers*, 1993.
- J. Pearl. On coding and filtering stationary signals by discrete fourier transforms (corresp.). *IEEE Transactions on Information Theory*, 19(2):229–232, 1973.
- F. Perronnin, J. Sánchez, and T. Mensink. Improving the Fisher kernel for large-scale image classification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2010.
- G. Peyré. Sparse modeling of textures. *Journal of Mathematical Imaging and Vision*, 34(1):17–31, 2009.
- D.-S. Pham and S. Venkatesh. Joint learning and dictionary construction for pattern recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008.
- D.-T. Pham. Fast algorithms for mutual information based independent component analysis. *IEEE Transactions on Signal Processing*, 52(10):2690–2700, 2004.
- N. Pinto, D. D. Cox, and J. J. DiCarlo. Why is real-world visual object recognition hard? *PLoS computational biology*, 4(1):151–156, 2008.

- J. Pokrass, A. M. Bronstein, M. M. Bronstein, P. Sprechmann, and G. Sapiro. Sparse modeling of intrinsic correspondences. In *Computer Graphics Forum*, volume 32, pages 459–468, 2013.
- M. Pontil, A. Argyriou, and T. Evgeniou. Multi-task feature learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2007.
- J. Portilla, V. Strela, M. J. Wainwright, and E. P. Simoncelli. Image denoising using scale mixtures of Gaussians in the wavelet domain. *IEEE Transactions on Image Processing*, 12(11):1338–1351, 2003.
- W. Pratt. Spatial transform coding of color images. *IEEE Transactions on Communication Technology*, 19(6):980–992, 1971.
- M. Protter and M. Elad. Image sequence denoising via sparse and redundant representations. *IEEE Transactions on Image Processing*, 18(1):27–35, 2009.
- H. Raguey, J. Fadili, and G. Peyré. A generalized forward-backward splitting. *SIAM Journal on Imaging Sciences*, 6(3):1199–1226, 2013.
- R. Raina, A. Battle, H. Lee, B. Packer, and A. Y. Ng. Self-taught learning: transfer learning from unlabeled data. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2007.
- I. Ramirez, F. Lecumberry, and G. Sapiro. Sparse modeling with universal priors and learned incoherent dictionaries. In *Proceedings of the 3rd IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, 2009.
- I. Ramirez, P. Sprechmann, and G. Sapiro. Classification and clustering via dictionary learning with structured incoherence and shared features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- M. Ranzato, F. J. Huang, Y-L. Boureau, and Y. LeCun. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007.
- G. Raskutti, M. J. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over-balls. *IEEE Transactions on Information Theory*, 57(10):6976–6994, 2011.
- B. Recht, C. Re, J. Tropp, and V. Bittorf. Factoring nonnegative matrices with linear programs. In *Advances in Neural Information Processing Systems (NIPS)*, 2012.

- X. Ren and D. Ramanan. Histograms of sparse codes for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- R. Rigamonti, M. A. Brown, and V. Lepetit. Are sparse representations really relevant for image classification? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2011.
- R. Rigamonti, A. Sironi, V. Lepetit, and P. Fua. Learning separable filters. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- J. Rissanen. Modeling by shortest data description. *Automatica*, 14(5):465–471, 1978.
- K. Ritter. Ein verfahren zur lösung parameterabhängiger, nichtlinearer maximum-probleme. *Mathematical Methods of Operations Research*, 6(4):149–166, 1962.
- F. Rodriguez and G. Sapiro. Sparse representations for image classification: Learning discriminative and reconstructive non-parametric dictionaries. Technical report, University of Minnesota, 2008.
- Y. Romano, M. Protter, and M. Elad. Single image interpolation via adaptive non-local sparsity-based modeling. *IEEE Transactions on Image Processing*, 23(7):3085–3098, 2014.
- S. Roth and M. J. Black. Fields of experts. *International Journal of Computer Vision*, 82(2):205–229, 2009.
- R. Rubinstein, M. Zibulevsky, and M. Elad. Efficient implementation of the K-SVD algorithm using batch orthogonal matching pursuit. Technical report, Technion - Computer Science Department, 2008.
- L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1-4):259–268, 1992.
- J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek. Image classification with the fisher vector: Theory and practice. *International Journal of Computer Vision*, 105(3):222–245, 2013.
- C. Schmid and R. Mohr. Local grayvalue invariants for image retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(5):530–535, May 1997.
- M. Schmidt, D. Kim, and S. Sra. Projected Newton-type methods in machine learning. In S. Sra, S. Nowozin, and S.J. Wright, editors, *Optimization for Machine Learning*. MIT Press, 2011.

- G. Schwarz. Estimating the dimension of a model. *Annals of Statistics*, 6(2): 461–464, 1978.
- M. W. Seeger. Bayesian inference and optimal design for the sparse linear model. *Journal of Machine Learning Research*, 9:759–813, 2008.
- T. Serre, L. Wolf, and T. Poggio. Object recognition with features inspired by visual cortex. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005.
- J. M. Shapiro. Embedded image coding using zerotrees of wavelet coefficients. *IEEE Transactions on Signal Processing*, 41(12):3445–3462, 1993.
- G. Sharma and H. J. Trussell. Digital color imaging. *IEEE Transactions on Image Processing*, 6(7):901–932, 1997.
- J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- E. P. Simoncelli and B. A. Olshausen. Natural image statistics and neural representation. *Annual review of neuroscience*, 24(1):1193–1216, 2001.
- E. P. Simoncelli, W. T. Freeman, E. H. Adelson, and D. J. Heeger. Shiftable multiscale transforms. *IEEE Transactions on Information Theory*, 38(2): 587–607, 1992.
- J. Sivic and A. Zisserman. Video Google: a text retrieval approach to object matching in videos. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2003.
- K. Skretting and K. Engan. Recursive least squares dictionary learning algorithm. *IEEE Transactions on Signal Processing*, 58(4):2121–2130, 2010.
- H. O. Song, S. Zickler, T. Althoff, R. Girshick, M. Fritz, C. Geyer, P. Felzenszwalb, and T. Darrell. Sparselet models for efficient multiclass object detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2012.
- D. A. Spielman, H. Wang, and J. Wright. Exact recovery of sparsely-used dictionaries. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2013.
- N. Srebro, J. D. M. Rennie, and T. S. Jaakkola. Maximum-margin matrix factorization. In *Advances in Neural Information Processing Systems (NIPS)*, 2005.
- J.-L. Starck, D. L. Donoho, and E. J. Candès. Astronomical image representation by the curvelet transform. *Astronomy and Astrophysics*, 398(2): 785–800, 2003.

- A. Szlam, M. Maggioni, and R.R. Coifman. Regularization on graphs with function-adapted diffusion processes. *Journal of Machine Learning Research*, 2007.
- H. Takeda, S. Farsiu, and P. Milanfar. Kernel regression for image processing and reconstruction. *IEEE Transactions on Image Processing*, 16(2):349–366, 2007.
- H. L. Taylor, S. C. Banks, and J. F. McCoy. Deconvolution with the ℓ_1 norm. *Geophysics*, 44(1):39–52, 1979.
- R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society Series B*, 58(1):267–288, 1996.
- R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused Lasso. *Journal of the Royal Statistical Society Series B*, 67(1):91–108, 2005.
- A. N. Tikhonov. Solution of incorrectly formulated problems and the regularization method. In *Soviet Math. Doklady.*, volume 4, pages 1035–1038, 1963.
- E. Tola, V. Lepetit, and P. Fua. DAISY: An efficient dense descriptor applied to wide-baseline stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(5):815–830, 2010.
- I. Tomic, B. A. Olshausen, and B. J. Culpepper. Learning sparse representations of depth. *IEEE Journal of Selected Topics in Signal Processing*, 5(5):941–952, 2011.
- J. A. Tropp. Greed is good: Algorithmic results for sparse approximation. *IEEE Transactions on Signal Processing*, 50(10):2231–2242, 2004.
- J. A. Tropp and A. C. Gilbert. Signal recovery from random measurements via orthogonal matching pursuit. *IEEE Transactions on Information Theory*, 53(12):4655–4666, 2007.
- J. A. Tropp, A. C. Gilbert, and M. J. Strauss. Algorithms for simultaneous sparse approximation. part I: Greedy pursuit. *Signal Processing, special issue "sparse approximations in signal and image processing"*, 86:572–588, 2006.
- P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*, 117(1):387–423, 2009.
- D. Y. Ts'o and C. D. Gilbert. The organization of chromatic and spatial interactions in the primate striate cortex. *The Journal of Neuroscience*, 8(5):1712–1727, 1988.

- A. B. Tsybakov. Optimal rates of aggregation. In *Proceedings of the Annual Conference on Computational Learning Theory*, 2003.
- B. A. Turlach, W. N. Venables, and S. J. Wright. Simultaneous variable selection. *Technometrics*, 47(3):349–363, 2005.
- D. Vainsencher, S. Mannor, and A. M. Bruckstein. The sample complexity of dictionary learning. *Journal of Machine Learning Research*, 12:3259–3281, 2011.
- J. C. van Gemert, C. J. Venman, A. W. M. Smeulders, and J. M. Geusebroek. Visual word ambiguity. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(7):1271–1283, 2010.
- G. Varoquaux, A. Gramfort, F. Pedregosa, V. Michel, and B. Thirion. Multi-subject dictionary learning to segment an atlas of brain spontaneous activity. In *Information Processing in Medical Imaging*, pages 562–573, 2011.
- G. Vinci, P. Freeman, J. Newman, L. Wasserman, and C. Genovese. Estimating the distribution of galaxy morphologies on a continuous space. In *Statistical Challenges in 21st Century Cosmology. Proceedings IAU Symposium No. 306*, 2014.
- H. von Helmholtz. LXXXI. on the theory of compound colours. *Philosophical Magazine Series 4*, 4(28):519–534, 1852.
- A. Wagner, J. Wright, A. Ganesh, Z. Zhou, H. Mobahi, and Y. Ma. Toward a practical face recognition system: robust alignment and illumination by sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(2):372–386, 2012.
- M.J. Wainwright. Sharp thresholds for noisy and high-dimensional recovery of sparsity using ℓ_1 -constrained quadratic programming. *IEEE Transactions on Information Theory*, 55(5):2183–2202, 2009.
- C. Wallraven, B. Caputo, and A. Graf. Recognition with local features: the kernel recipe. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2003.
- J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong. Locality-constrained linear coding for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- S. Wang, D. Zhang, Y. Liang, and Q. Pan. Semi-coupled dictionary learning with applications to image super-resolution and photo-sketch synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

- X. Wang and X. Tang. Face photo-sketch synthesis and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(11):1955–1967, 2009.
- L. Wasserman. *All of nonparametric statistics*. Springer, 2006.
- G. S. Watson. Smooth regression analysis. *Sankhya, The Indian Journal of Statistics, Series A*, 26:359–372, 1964.
- J. Wright, A.Y. Yang, A. Ganesh, S.S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2):210–227, 2009.
- D. Wrinch and H. Jeffreys. XLII. On certain fundamental principles of scientific inquiry. *Philosophical Magazine Series 6*, 42(249):369–390, 1921.
- T. T. Wu and K. Lange. Coordinate descent algorithms for lasso penalized regression. *Annals of Applied Statistics*, 2(1):224–244, 2008.
- Y. Wu, J. Lim, and M.-H. Yang. Online object tracking: A benchmark. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 9:2543–2596, 2010.
- R. Xiaofeng and L. Bo. Discriminatively trained sparse code gradients for contour detection. In *Advances in Neural Information Processing Systems (NIPS)*, 2012.
- Z. Xing, M. Zhou, A. Castrodad, G. Sapiro, and L. Carin. Dictionary learning for noisy and incomplete hyperspectral images. *SIAM Journal on Imaging Sciences*, 5(1):33–56, 2012.
- J. Yang, K. Yu, Y. Gong, and T. Huang. Linear spatial pyramid matching using sparse coding for image classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 19(11):2861–2873, 2010a.
- J. Yang, K. Yu, and T. Huang. Supervised translation-invariant sparse coding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010b.
- J. Yang, Z. Wang, Z. Lin, S. Cohen, and T. Huang. Coupled dictionary training for image super-resolution. *IEEE Transactions on Image Processing*, 21(8):3467–3478, 2012a.

- J. Yang, K. Yu, and T. Huang. Efficient highly overcomplete sparse coding using a mixture model. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2012b.
- M. Yang, D. Zhang, and X. Feng. Fisher discrimination dictionary learning for sparse representation. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011.
- T. Young. *A course of lectures on natural philosophy and the mechanical art*. Taylor and Watson, 1845.
- G. Yu, G. Sapiro, and S. Mallat. Solving inverse problems with piecewise linear estimators: from Gaussian mixture models to structured sparsity. *IEEE Transactions on Image Processing*, 21(5):2481–2499, 2012.
- K. Yu and T. Zhang. Improved local coordinate coding using local tangents. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2010.
- K. Yu, T. Zhang, and Y. Gong. Nonlinear learning using local coordinate coding. In *Advances in Neural Information Processing Systems (NIPS)*, 2009.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society Series B*, 68: 49–67, 2006.
- M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014.
- M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus. Deconvolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- M. D. Zeiler, G. W. Taylor, and R. Fergus. Adaptive deconvolutional networks for mid and high level feature learning. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011.
- R. Zeyde, M. Elad, and M. Protter. On single image scale-up using sparse representations. In *Curves and Surfaces*, pages 711–730. Springer, 2012.
- Q. Zhang and B. Li. Discriminative K-SVD for dictionary learning in face recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- Y. Zhang, M. J. Wainwright, and M. I. Jordan. Lower bounds on the performance of polynomial-time algorithms for sparse linear regression. *preprint arXiv:1402.1918*, 2014.

- 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.
- M. Zhou, H. Chen, L. Ren, G. Sapiro, L. Carin, and J. W. Paisley. Nonparametric bayesian dictionary learning for sparse image representations. In *Advances in Neural Information Processing Systems (NIPS)*, 2009.
- M. Zhou, H. Chen, J. Paisley, L. Ren, L. Li, Z. Xing, D. Dunson, G. Sapiro, and L. Carin. Nonparametric Bayesian dictionary learning for analysis of noisy and incomplete images. *IEEE Transactions on Image Processing*, 21(1):130–144, 2012.
- X. Zhou, K. Yu, T. Zhang, and T. Huang. Image classification using super-vector coding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2010.
- S. C. Zhu and D. Mumford. Prior learning and gibbs reaction-diffusion. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(11):1236–1250, 1997.
- S.-C. Zhu, C.-E. Guo, Y. Wang, and Z. Xu. What are textons? *International Journal of Computer Vision*, 62(1-2):121–143, 2005.
- D. Zoran and Y. Weiss. From learning models of natural image patches to whole image restoration. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011.
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B*, 67(2):301–320, 2005.