

Global Optimality of Sparse Dictionary Learning with Applications to Subspace Clustering

by

Imke Mayer

A thesis submitted to École Normale Supérieure Paris Saclay in partial
fulfillment of the requirements for the degree of Master of Science.

Master's program "Mathématiques Vision Apprentissage"

Center for Imaging Science
Johns Hopkins University, Baltimore

September 2018

école —————
normale —————
supérieure —————
paris-saclay —————

Abstract

Most signal processing applications require a good choice of representation of the data for numerous reasons. For instance, classical tasks such as compression, denoising or classification can be performed more easily when the data is represented in terms of certain bases that exploit special structure in high-dimensional data. One specific representation is obtained by decomposing the data in terms of a dictionary, whose elements are either predefined such as wavelets or learned directly from the data. When the representation is assumed to be sparse, i.e., when each data point is a linear combination of a few elements from the dictionary, the latter problem is called Sparse Dictionary Learning (SDL). This class of data modeling methods exploits certain structure in the data which can be formulated as: *the data lies in a union of low-dimensional subspaces* and therefore the dictionary elements, called *atoms*, can be seen as basis elements of these subspaces. Since the dictionary is typically overcomplete, the atoms are actually dependent and should be considered as frames for the subspaces that capture the main directions in terms of compact representations. While a variety of efficient algorithms that give approximate solutions for the SDL problem have been proposed, the theoretical analysis of this problem remains a key challenge due to its non-convexity. In this work we present a step towards a theoretical understanding of global optimality for the SDL problem. In particular, we show that by regularizing the SDL problem with a certain function that penalizes the number of dictionary atoms it is possible to derive conditions under which an optimal dictionary can be found. We also demonstrate this analysis with an application of SDL to the problem of Subspace Clustering.

Primary Reader and Advisor: Prof. René Vidal

Secondary Readers: Dr. Benjamin Haeffele, Dr. Benjamín Béjar Haro

Acknowledgments

This thesis would not have been possible without the guidance of my supervisors Prof. René Vidal and Dr. Benjamin Haeffele. I would like to thank you for giving me the opportunity to contribute to your project.

Furthermore, I would like to say thank you to Connor Lane who provided me with valuable advices on theoretic as well as more technical aspects of this work.

That said, I would like to give thanks to the entire Vision, Dynamics and Learning Lab for giving me the chance to gain insights into their research projects and into their day-to-day work as researchers in a multidisciplinary environment.

Contents

Abstract	ii
Acknowledgments	iii
1 Introduction	1
1.1 Problem formulation	2
1.2 Related work	5
1.3 Contributions	7
2 Global Optimality for Sparse Dictionary Learning	9
2.1 Review of Conditions for Global Optimality in Structured Matrix Factorization	9
2.2 Structured matrix factorization algorithm for SDL	15
2.2.1 Meta-algorithm	16
2.2.2 Details on local optimization strategies	17
2.3 Analysis of the polar problem for SDL	25
2.3.1 Empirical estimation of the optimization geometry	27
2.3.2 Restricted necessary optimality conditions	28
2.4 Extension to discriminative sparse dictionary learning	32
3 Subspace Clustering and Sparse Dictionary Learning	36
3.1 Subspace clustering by sparse or low-rank representation	36
3.1.1 Sparse subspace clustering	38
3.1.2 Low-rank subspace clustering	39
3.2 Subspace clustering by sparse dictionary learning	40
3.3 Experiments	44
3.3.1 Preliminaries	44

CONTENTS

3.3.2 Recovering the generating factorization size	47
3.3.3 Subspace clustering	52
4 Conclusion	55
A Appendix	56
Bibliography	60

1 Introduction

Modeling data as sparse linear combinations of some basis elements is a widely used approach to deal with complex data and is motivated by different applications in various domains. For example in signal processing, sparse representations can allow for better results in signal denoising and compression [10, 27], machine learning sparse representations can lead to better results in classification [26, 37], and in statistics sparse modeling can be used for variable selection resulting in more interpretable statistical models [35]. The main assumption and goal of such models can be described as follows: the relevant information contained in the (often high-dimensional) data often lies in a low-dimensional structure and finding this low-dimensional structure can therefore allow for more efficient, robust and interpretable data processing. If one assumes that the data can be modeled by a linear low-dimensional structure, i.e. by one or more linear subspaces, the model which allows recovering such structure is called *structured matrix factorization*. In this section, we give the formal definition of structured matrix factorization, some of its well-known applications and we cast a specific class of sparse models called *sparse dictionary learning* as a special case of structured matrix factorization. In sparse dictionary learning each data point is decomposed as a sparse linear combination of elements from a dictionary, and both the dictionary and the sparse representations are learned from the data. We will discuss its applicability to another special case of structured matrix factorization, the problem of *subspace clustering*.

1.1 Problem formulation

Let $X \in \mathbb{R}^{D \times N}$ be the data matrix which can be considered as the concatenation of N vectorized signals $\{x_i, \dots, x_N\} \subset \mathbb{R}^D$. Under the linear data model, i.e. the assumption that the data can be approximated by one or more linear subspaces, structured matrix factorization consists in decomposing the data X into two factors U and V in a way that translates or exploits the low-dimensionality assumptions on the data. More formally, given a loss function ℓ which measures the quality of approximation of the factorization UV^T w.r.t. the initial data X and a regularization function Θ which enforces specific structure on the factors U and V , we can write the structured matrix factorization problem as:

$$\min_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r}}} \ell(X, UV^T) + \lambda \Theta(U, V), \quad (1.1)$$

where we usually assume $r < \min\{D, N\}$. The interest of structured matrix factorization is two-fold: especially in the case of high-dimensional data it can considerably reduce the representation complexity of the data (i.e. $\mathcal{O}((D + N)r)$ instead of $\mathcal{O}(DN)$). Furthermore, it can allow to separate the information contained in the data into the different factors by enforcing certain constraints on each factor, such as spatial and temporal information in the case of videos as demonstrated in [18] and therefore reveal interesting patterns in the data. The generic form of (1.1) allows to derive other possible applications and to cast other well-known models as special cases of this problem.

We will now turn to the problem of sparse dictionary learning and explain the choices of loss and regularization function for (1.1) in this case. Given a set of N signals $\{x_1, \dots, x_N\} \subset \mathbb{R}^D$, the dictionary learning problem consists of finding a set of dictionary atoms $\{d_1, \dots, d_r\} \subset \mathbb{R}^D$ and N sparse codes $\{a_1, \dots, a_N\} \subset \mathbb{R}^r$ that allow to recover, either exactly or approximately, the initial set of signals (which can be one-dimensional but also vectorized forms of multi-dimensional data such as images or videos). In other words, it requires solving the following problem:

CHAPTER 1. INTRODUCTION

$$\min_{\substack{\mathbf{d}_1, \dots, \mathbf{d}_r \in \mathbb{R}^D \\ \mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^r}} \frac{1}{2} \sum_{i=1}^N \left(\mathbf{x}_i - \sum_{j=1}^r (\mathbf{a}_i)_j \mathbf{d}_j \right)^2 + \lambda \Theta(\{\mathbf{d}_k\}_{k=1}^r, \{\mathbf{a}_k\}_{k=1}^r), \quad (1.2)$$

where Θ is defined on the dictionary atoms and the sparse codes and $\lambda \geq 0$ is a parameter adjusting the trade-off between data fidelity and regularization.

In matrix form this problem can be rewritten as

$$\min_{\substack{D \in \mathbb{R}^{D \times r} \\ A \in \mathbb{R}^{r \times N}}} \frac{1}{2} \|X - DA\|_F^2 + \lambda \Theta(D, A), \quad (1.3)$$

where the r columns of D are the dictionary atoms and the N columns of A contain the sparse representation of the data $X = [\mathbf{x}_1 \dots \mathbf{x}_N] \in \mathbb{R}^{D \times N}$ in terms of the dictionary D .

Without any regularization, i.e. if $\lambda = 0$, this problem can be an ill-posed signal reconstruction problem. For instance, given a signal $\mathbf{x} \in \mathbb{R}^D$ and a dictionary D with non-trivial null-space, there are infinitely many solutions $\mathbf{a} \in \mathbb{R}^r$ to the problem $\min_{\mathbf{a}} \frac{1}{2} \|\mathbf{x} - D\mathbf{a}\|_2^2$. But assuming that the decomposition of the data into D and A is such that each signal only “uses” very few atoms, i.e. the columns of A are sparse, and that the dictionary is overcomplete ($r > D$), the problem becomes the *Sparse Dictionary Learning* problem (SDL), first introduced in [28] to efficiently process natural image patches. Hence the regularization function Θ can be considered to be defined on the $\ell_{1,1}$ norm of matrix A , i.e. on the sum of the absolute values of all its entries, since it is well known that the ℓ_1 norm is a sparsity-inducing convex relaxation of the pseudo-norm ℓ_0 . The resulting decomposition of the data into dictionary and codes is illustrated in Figure 1.1.

Another well-known class of problems which exploits (and recovers) a certain structure of the – often high-dimensional – data and which can be cast as a matrix factorization model is Subspace Clustering (SC). For this class of problems the goal is to recover a low-dimensional structure which the data comes

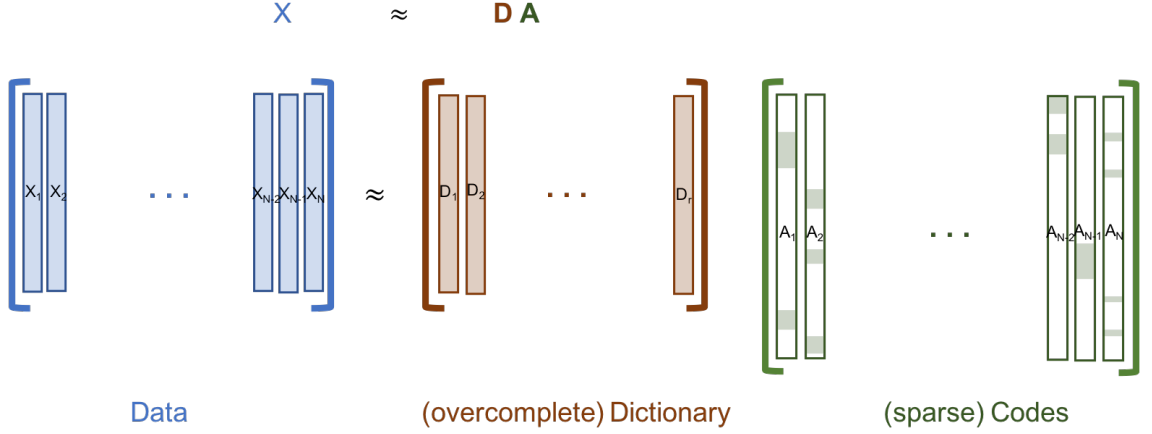


Figure 1.1: Sparse dictionary learning model.

from, more specifically a union of multiple low-dimensional subspaces, allowing for a more robust and compact representation of the data. Since this is a combinatorial problem which is NP-hard, current subspace clustering methods (see [36] for a review) only seek to find a possible representation of the subspaces (a large dictionary with linearly dependent atoms) and to identify which subspace a data point comes from, based on such a dictionary. The subspaces' representation being redundant, one cannot hope for a unique sparse representation.

While the applications of structured matrix factorization, and of sparse dictionary learning and subspace clustering as special cases, are numerous, the main issue is that it is a non-convex problem due to the matrix multiplication in the data fidelity term of the objective function. Therefore most optimization strategies require good heuristics for the choice of initialization and there are no guarantees on global optimality. But if we keep the dictionary D in (1.3) fixed and assume that Θ is convex w.r.t. A , the problem reduces to a regularized least squares problem in A and is therefore a convex problem (for specific choices of Θ such as the ℓ_1 norm numerous efficient algorithms solving this lasso-type problem, known as *sparse coding*, have been proposed [9, 25]). However, analyzing and solving the joint optimization problem is much more challenging and is still an open field of research. Motivated by the work of [4] and

CHAPTER 1. INTRODUCTION

more specifically [18] who proposed a general matrix factorization framework to characterize global optimality and several explicit methods to find globally optimal solutions in certain cases of structured matrix factorization, we follow the same line of approach to analyze the non-convex SDL problem.

1.2 Related work

MATRIX FACTORIZATION AND SPARSE DICTIONARY LEARNING

The analysis of the optimization landscape of non-convex decomposition or signal recovery problems has recently gained much attention due to the recurring empirical observation that many apparently hard problems can be reasonably well solved by approximate local strategies which do not come with a guarantee to converge to a “good” solution of non-convex problems but only to some local optima or critical points. This suggests that the optimization landscape of these problems is favorable in the sense that there are no spurious local minima or bad saddle points and therefore any descent strategy can achieve a good result without getting stuck in a bad critical point.

For the case of non-convex low-rank matrix optimization problems [38] showed that there are no spurious local minima and extend the concept of strict saddle property introduced by [13] implying that most optimization strategies such as gradient descent converge to a global solution from any initialization. The problem of optimal dictionary recovery has been studied by [34] in the complete case ($r = D$) over the sphere and they propose an algorithm that provably recovers an optimal sparsifying dictionary, again by using the argument that there are no spurious local minima. Extending the ideas of [17,18], the authors of [32] propose a method with guarantees on global optimality to solve the separable dictionary learning problem applied to tensors. Another approach to (approximately) solve such non-convex matrix factorization problems is to introduce some convex relaxation (such as the nuclear norm in the low-rank matrix factorization case, [6]) to transform the factorization problem into a convex

CHAPTER 1. INTRODUCTION

matrix approximation problem. Convex relaxations are an appealing solution from the methodological point of view since they allow an easier analysis of problem and its optimal solutions and can often be solved by using standard optimization tools (efficient MATLAB implementations are provided for instance by [15]). A caveat however is that these formulations often do not scale to larger dataset, i.e. their application to large scale datasets is challenging and often inadvisable from a practical point of view. Furthermore the use of (tight) convex relaxations cannot capture the entire modeling flexibility of the structured matrix factorization formulation. For instance, certain structures on the factors U and V , such as non-negativity [22], cannot be promoted by regularizing their product UV^T .

SUBSPACE CLUSTERING

Many different approaches have been proposed for solving the subspace clustering problem, ranging from a generalization of the well-known k-means to the multiple subspace setting, to algebraic, statistical and spectral methods (see [36] for a review). Among the spectral methods the Sparse Subspace Clustering (SSC, [11]) is particularly efficient (it only requires to solve a convex ℓ_1 minimization problem) and robust to different types of contaminated data such as noise, missing data and sparse outliers. The problem of sparse subspace clustering is also non-convex but under certain assumptions on the data distribution and subspace configuration the solution to the convex surrogate problem induced by the ℓ_1 norm provides the desired sparse representation of the data. Another assumption which can be made is that the data allows a low-rank representation [23] and in the case of noise-free data drawn from independent subspaces, the optimal solution to the corresponding optimization problem can be written in closed form from the SVD of the data. These methods consider the data as a self-expressive dictionary, the factorization of the initial data X into two factors therefore does not come with any (memory) complexity reduction. [3] however propose a subspace clustering method which is

CHAPTER 1. INTRODUCTION

based on a smaller dictionary and the clustering is then made possible through bipartite graph modeling.

1.3 Contributions

In this work we present a theoretical and experimental study of the problem of structured matrix factorization in the case of sparse dictionary learning. We focus on the problem of solving an auxiliary optimization problem, the *polar problem*, which serves as a certificate of global optimality for the initial SDL problem. We give some intuitions and first steps towards a theoretical analysis of the polar problem for the sparse dictionary learning problem and discuss the impact of approximation errors on the initial problem. We also experimentally demonstrate its relevance for solving the overall matrix factorization formulation, especially for recovering the “generating” factorization size on data drawn from a union of subspaces.

In Chapter 2 we provide the theoretical analysis of global optimality for the SDL problem, discuss and analyze the challenges of verifying the derived global optimality conditions and propose an iterative algorithm for SDL based on the previous analysis. In Chapter 3 we draw a connection of our SDL model to existing subspace clustering models to motivate a novel approach to subspace clustering via our SDL matrix factorization formulation.

Notations

In the next chapters, we will use the following fairly standard notations: $[n]$ denotes the set of positive integers up to n , $\{1, \dots, n\}$; for a set $S \subseteq \Omega$ we denote its complement by $S^c = \Omega \setminus S \subseteq \Omega$; scalars are denoted by regular lower case letters, vectors by bold lower case letters; matrices are denoted by capital letters. Columns and rows of a matrix $X \in \mathbb{R}^{m \times n}$ are referenced by X_j with $j \in [n]$ and $X_{<i>}$ with $i \in [m]$ respectively. For a scalar a , we define $(a)_+ := \max(a, 0)$.

CHAPTER 1. INTRODUCTION

Unless specified differently, all vectors are considered as column vectors. The notation $\|\cdot\|$ will be used indifferently for both vector and matrix norms. For instance $\|\cdot\|_F$ denotes the Frobenius norm on matrices, $\|\cdot\|_{k \rightarrow l}$ the $k \rightarrow l$ operator norm and $\|\cdot\|_{k,l}$ the $L_{k,l}$ norm ($k, l \in \mathbb{N}^*$).

2 Global Optimality for Sparse Dictionary Learning

In this chapter we will analyze the problem of Sparse Dictionary Learning as a particular case of regularized matrix factorization problem and characterize its global optima. Based on these characterizations of global optimality we propose an iterative algorithm for solving the non-convex optimization problem. We will focus on an auxiliary optimization problem which arises in these characterizations of global optima. We empirically assess the difficulty of solving this problem, propose a conjecture on its optimization landscape and present initialization and optimization strategies which empirically lead to good approximate solutions.

2.1 Review of Conditions for Global Optimality in Structured Matrix Factorization

We recall the standard formulation of sparse dictionary learning which we want to analyze as a special case of the structured matrix factorization problem.

$$\min_{\substack{U \in \mathbb{R}^{D \times r} \\ V^T \in \mathbb{R}^{r \times N}}} \frac{1}{2} \|X - UV^T\|_F^2 + \lambda \|V\|_{1,1} \quad \text{subject to } \|U_i\|_2 \leq 1 \text{ for all } i \in [r]. \quad (2.1)$$

The objective is a sum of a differentiable loss function and a convex but non-differentiable regularization function. The loss function is convex w.r.t. the

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

product UV^T and w.r.t. each factor but it is not jointly convex in (U, V) , making the overall optimization problem non-convex. Minimization is defined over U and V which are of fixed size, i.e. the size r of the factorization is part of the problem dimensions $(D, N, r) \in (\mathbb{N}_+)^3$.

Before we derive our sparse dictionary learning model from the above formulation and delve into its analysis we begin by defining some relevant notations and quantities, following the work of [17, 18].

Definition 1 (Rank-1 regularizer, [18]). A function $\theta : \mathbb{R}^D \times \mathbb{R}^N \rightarrow \mathbb{R}_+ \cup \infty$ is said to be a rank-1 regularizer if

1. $\theta(\mathbf{u}, \mathbf{v})$ is positively homogeneous with degree $p > 0$, i.e. $\theta(\alpha \mathbf{u}, \alpha \mathbf{v}) = \alpha^p \theta(\mathbf{u}, \mathbf{v})$, $\forall \alpha \geq 0, \forall (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^D \times \mathbb{R}^N$.
2. $\theta(\mathbf{u}, \mathbf{v})$ is positive semi-definite, i.e. $\theta(\mathbf{0}, \mathbf{0}) = 0$ and $\theta(\mathbf{u}, \mathbf{v}) \geq 0$ for all (\mathbf{u}, \mathbf{v}) .
3. For any sequence $(\mathbf{u}_n, \mathbf{v}_n)$ such that $\|\mathbf{u}_n \mathbf{v}_n^T\| \rightarrow \infty$ we have that $\theta(\mathbf{u}_n, \mathbf{v}_n) \rightarrow \infty$.

Definition 2 (Elemental mapping, r -element factorization mapping, [18]). An elemental mapping $\phi : \mathbb{R}^D \times \mathbb{R}^N \rightarrow \mathbb{R}^{D \times N}$ is any mapping which is positively homogeneous with degree $p > 0$. The r -element factorization mapping $\Phi_r : \mathbb{R}^{D \times r} \times \mathbb{R}^{N \times r} \rightarrow \mathbb{R}^{D \times N}$ is defined as $\Phi_r(U, V) = \sum_{i=1}^r \phi(U_i, V_i)$.

Unless specified differently we will use $\phi : (\mathbf{u}, \mathbf{v}) \mapsto \mathbf{u} \mathbf{v}^T$ and $\Phi_r : (U, V) \mapsto UV^T = \sum_i U_i V_i^T$ as elemental mapping and r -element factorization mapping respectively.

These definitions can be extended to higher dimensions, for instance to analyze more general tensor decomposition problems as detailed in [18], but here we limit the presentation to the matrix factorization case as it is sufficient to analyze the problem of sparse dictionary learning.

Definition 3 (Matrix factorization regularizer, [18]). Given a rank-1 regularizer θ which is positively homogeneous with degree 2, the matrix factorization

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

regularizer $\Omega_\theta : \mathbb{R}^{D \times N} \rightarrow \mathbb{R}_+ \cup \infty$ is defined as

$$\Omega_\theta(Y) = \inf_{r \in \mathbb{N}} \inf_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r}}} \sum_{i=1}^r \theta(U_i, V_i) \quad \text{subject to } Y = UV^T. \quad (2.2)$$

An important remark on this last definition is that such a matrix factorization regularizer will be convex w.r.t. Y . And we can also note that if θ is such that $\theta(-\mathbf{u}, \mathbf{v}) = \theta(\mathbf{u}, \mathbf{v})$ or $\theta(\mathbf{u}, -\mathbf{v}) = \theta(\mathbf{u}, \mathbf{v})$ for all (\mathbf{u}, \mathbf{v}) , then Ω_θ is a norm on X , called *atomic norm* or *decomposition norm*. For instance, taking $\theta(\mathbf{u}, \mathbf{v}) = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2$, this definition corresponds to the variational form of the nuclear norm $\|\cdot\|_*$ for matrices which is defined as the sum of the singular values of a matrix.

These definitions allow to introduce the framework proposed in [18] with the following key idea: the matrix factorization regularizer allows to tightly couple the structured matrix factorization problem which is non-convex in (U, V) to a problem that is convex in Y . Exploiting the convexity of the latter allows to analyze the former and to derive conditions of globally optimality for it. An important aspect which allows this coupling is to consider the factorization size r as one of the variables and therefore to optimize over factorizations of all possible sizes $r \in \mathbb{N}_+$.

We can now go on to analyze how the sparse dictionary learning problem can be cast into this matrix factorization framework in order to give global optimality conditions for the non-convex sparse dictionary learning problem. As pointed out in [5], constraining the dictionary atoms in problem (2.1) to have ℓ_2 -norm less than 1 is equivalent to penalizing their ℓ_2 -norm by integrating these terms in the regularization term of the objective function. Indeed, this becomes clear if we write the variational form of the $\ell_{2,1}$ norm of some matrix $M = [M_1 \dots M_N]$

$$\|M\|_{2,1} = \sum_{i=1}^N \|M_i\|_2 = \frac{1}{2} \min_{\eta_i \geq 0} \sum_{i=1}^N \frac{\|M_i\|_2^2}{\eta_i} + \eta_i, \quad (2.3)$$

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

which is obtained by a simple application of the Cauchy-Schwarz inequality. Additionally we will allow the factorization size r to vary, since we argued earlier that this will allow to couple this non-convex factorization problem to a closely related convex problem. This leads to a first formulation:

$$\min_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r} \\ r \in \mathbb{N}_+}} \frac{1}{2} \|X - UV^T\|_F^2 + \lambda \sum_{i=1}^r \|U_i\|_2 \|V_i\|_1 \quad (2.4)$$

Observe that the factorization of a matrix Y as the product UV^T is not unique due to a rotational ambiguity, i.e. if $(U, V) \in \mathbb{R}^{D \times r} \times \mathbb{R}^{N \times r}$ are such that $Y = UV^T$, then for any orthogonal matrix $R \in O(r)$, we have $Y = URR^TV^T$. The regularizer can eliminate this ambiguity by enforcing special structure on the factors U and V , which in turn enforces special structure on the column and row spaces of UV^T , respectively. For instance, using the ℓ_1 norm on the columns of V , i.e. in the row space, the decomposition of the data X over the dictionary will be very sparse and the dictionary will consequently be large since we optimize over all possible factorization sizes. Conversely, using the ℓ_2 norm on the columns of V will lead to a compact dictionary but the codes (rows of V) will not be sparse. In order to define a trade-off between these two boundary cases we modify the regularization function and propose to solve the alternative problem:

$$\min_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r} \\ r \in \mathbb{N}_+}} \underbrace{\frac{1}{2} \|X - UV^T\|_F^2}_{\ell(X, \Phi_r(U, V))} + \lambda \underbrace{\sum_{i=1}^r \|U_i\|_2 (\gamma \|V_i\|_1 + (1 - \gamma) \|V_i\|_2)}_{\Theta_\gamma(U, V)} \quad (2.5)$$

Again we draw the attention to the fact that in the above formulation minimization is carried out over the factors and the factorization size r . This implies that the optimality conditions we will present next are also defined on the factorization size, superseding the model selection step that requires the use of some heuristics in most applications to define a fixed size of the factors prior to solving the factorization problem. Instead, the size of the factorization

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

is data-driven through the use of the regularization function Θ_γ .

The loss term of the objective function (2.5) remains unchanged from the previous formulations (2.1) and (2.4). But we introduce a new regularization function that depends on the parameter $\gamma \in [0, 1]$. Observe that the new regularization function is defined as the sum of a rank-1 regularizer over the columns of U and V . Indeed we easily see that

$$\theta_\gamma : (\mathbf{u}, \mathbf{v}) \in \mathbb{R}^D \times \mathbb{R}^N \mapsto \|\mathbf{u}\|_2(\gamma\|\mathbf{v}\|_1 + (1 - \gamma)\|\mathbf{v}\|_2) \in \mathbb{R}_+$$

is positively homogeneous with degree 2, positive semi-definite and such that for any sequence $(\mathbf{u}_n, \mathbf{v}_n)$ such that $\|\mathbf{u}_n \mathbf{v}_n^T\| \rightarrow \infty$, we have $\theta_\gamma(\mathbf{u}_n, \mathbf{v}_n) \rightarrow \infty$ (since for all $n \in \mathbb{N}$, $\|\mathbf{u}_n \mathbf{v}_n^T\| = \|\mathbf{u}_n\|_2 \|\mathbf{v}_n\|_2 \leq \|\mathbf{u}_n\|_2 \|\mathbf{v}_n\|_1$ and therefore $\|\mathbf{u}_n \mathbf{v}_n^T\| \leq \theta_\gamma(\mathbf{u}_n, \mathbf{v}_n)$).

We remark that the regularization function Θ_γ is closely related to the matrix factorization regularizer (2.2), also referred to as decomposition norm or atomic norm in the sense that we always have

$$\underbrace{\ell(X, Y) + \lambda \Omega_{\theta_\gamma}(Y)}_{=: F(Y)} \leq \underbrace{\ell(X, \Phi_r(U, V)) + \lambda \Theta_\gamma(U, V)}_{=: f(U, V)} \quad (2.6)$$

$\forall (U, V)$ such that $Y = \Phi_r(U, V) = UV^T$.

The main result of [18] makes use of this global lower bound of the non-convex objective function f (in the factor space) by the convex function F (in the product space) to prove the following results:

Theorem 1 (Theorem 1 in [18]). *Given a function $\ell(X, Y)$ that is convex in Y and once differentiable w.r.t. Y , a rank-1 regularizer θ_γ and a constant $\lambda > 0$, local minima (\tilde{U}, \tilde{V}) of (2.5) are globally optimal if $(\tilde{U}_i, \tilde{V}_i) = (0, 0)$ for some $i \in [r]$. Moreover, $\hat{Y} = \Phi_r(\tilde{U}, \tilde{V}) = \tilde{U}\tilde{V}^T$ is global minima of $F(Y)$ and $\tilde{U}\tilde{V}^T$ is an optimal factorization of \hat{Y} , i.e. it achieves the infimum in the definition (2.2) of the atomic norm $\Omega_{\theta_\gamma}(\hat{Y})$.*

This theorem roughly states that under the given assumptions any local

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

minima of f which is sufficiently large is a global minima. And the following corollary suggests a way to verify that a given point is a local minima and sufficiently large, hence a global minima.

Corollary 1 (Corollary 1 in [18]). *Under the same assumptions as in the previous theorem, a local minima (\tilde{U}, \tilde{V}) of $f(U, V)$ is globally optimal if it satisfies:*

1. $\tilde{U}_i^T \left(-\frac{1}{\lambda} \nabla_{\Phi_r} \ell(X, \Phi_r(\tilde{U}, \tilde{V})) \right) \tilde{V}_i = \theta_\gamma(\tilde{U}_i, \tilde{V}_i) \quad \text{for all } i = 1, \dots, r$
2. $\mathbf{u}^T \left(-\frac{1}{\lambda} \nabla_{\Phi_r} \ell(X, \Phi_r(\tilde{U}, \tilde{V})) \right) \mathbf{v} \leq \theta_\gamma(\mathbf{u}, \mathbf{v}) \quad \text{for all } (\mathbf{u}, \mathbf{v}).$

It can be shown that for our choice of rank-1 regularizer, θ_γ , any first-order optimal point of (2.5) satisfies condition 1 [18, Proposition 3]¹. The second condition implies that at (\tilde{U}, \tilde{V}) , the addition of any new direction (\mathbf{u}, \mathbf{v}) , i.e. $(U, V) \leftarrow ([\tilde{U} \sqrt{\alpha} \mathbf{u}], [\tilde{V} \sqrt{\alpha} \mathbf{v}])$ for some $\alpha > 0$, will only increase the objective function f . This can be seen more easily by rearranging the terms of the inequality above:

$$\mathbf{u}^T \left(-\frac{1}{\lambda} \nabla_{\Phi_r} \ell(X, \Phi_r(\tilde{U}, \tilde{V})) \right) \mathbf{v} \leq \theta_\gamma(\mathbf{u}, \mathbf{v}) \Leftrightarrow 0 \leq \langle \nabla_{\Phi_r} \ell(X, \Phi_r(\tilde{U}, \tilde{V})), \mathbf{u} \mathbf{v}^T \rangle + \lambda \theta_\gamma(\mathbf{u}, \mathbf{v}),$$

where the right hand-side of the second equivalent term corresponds to the directional derivative of the objective function in the direction (\mathbf{u}, \mathbf{v}) . Therefore if condition 2 of the corollary holds we have $f([\tilde{U} \sqrt{\alpha} \mathbf{u}], [\tilde{V} \sqrt{\alpha} \mathbf{v}]) \geq f(\tilde{U}, \tilde{V})$, since by using the convexity of the loss ℓ w.r.t. to its second argument we have:

$$\begin{aligned} f([\tilde{U} \sqrt{\alpha} \mathbf{u}], [\tilde{V} \sqrt{\alpha} \mathbf{v}]) &= \ell(X, \tilde{U} \tilde{V}^T + \alpha \mathbf{u} \mathbf{v}^T) + \lambda \sum_{i=1}^r \theta(\tilde{U}_i, \tilde{V}_i) + \alpha \theta(\mathbf{u}, \mathbf{v}) \\ &\geq \ell(X, \tilde{U} \tilde{V}^T) + \alpha \langle \nabla_{\Phi_r} \ell(X, \Phi_r(\tilde{U}, \tilde{V})), \mathbf{u} \mathbf{v}^T \rangle \\ &\quad + \lambda \sum_{i=1}^r \theta_\gamma(\tilde{U}_i, \tilde{V}_i) + \alpha \lambda \theta_\gamma(\mathbf{u}, \mathbf{v}) \\ &\geq f(\tilde{U}, \tilde{V}). \end{aligned}$$

¹This proposition actually covers a broader class of possible rank-1 regularizers that encompasses the product of norms as a special case.

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

If condition 2 is violated, i.e. we find a pair (\mathbf{u}, \mathbf{v}) that does not satisfy the inequality in Corollary 1, we have found a descent direction that allows to decrease the objective by exactly $\alpha \left(\langle \nabla_{\Phi_r} \ell(X, \Phi_r(\tilde{U}, \tilde{V})), \mathbf{u}\mathbf{v}^T \rangle + \lambda \theta_\gamma(\mathbf{u}, \mathbf{v}) \right)$ where α is some (small) positive step-size.

The authors of [18] point out that condition 2 of the corollary can be verified by evaluating the so-called *polar function* $\Omega_{\theta_\gamma}^\circ$ at $Z = -\frac{1}{\lambda} \nabla_{\Phi_r} \ell(X, \Phi_r(\tilde{U}, \tilde{V}))$ and testing whether $\Omega_{\theta_\gamma}^\circ(Z) \leq 1$ where $\Omega_{\theta_\gamma}^\circ(Z)$ is defined as

$$\Omega_{\theta_\gamma}^\circ(Z) = \sup_{\mathbf{u}, \mathbf{v}} \mathbf{u}^T Z \mathbf{v} \quad \text{s.t. } \theta_\gamma(\mathbf{u}, \mathbf{v}) \leq 1. \quad (2.7)$$

Indeed, this becomes clear when writing out the following equivalences:

$$\mathbf{u}^T Z \mathbf{v} \leq \theta_\gamma(\mathbf{u}, \mathbf{v}) \quad \text{for all } (\mathbf{u}, \mathbf{v}) \quad (2.8)$$

$$\Leftrightarrow \frac{\mathbf{u}^T Z \mathbf{v}}{\theta_\gamma(\mathbf{u}, \mathbf{v})} \leq 1 \quad \text{for all } (\mathbf{u}, \mathbf{v}) \quad (2.9)$$

$$\Leftrightarrow \sup_{\substack{(\mathbf{u}, \mathbf{v}) \\ \theta_\gamma(\mathbf{u}, \mathbf{v}) \leq 1}} \mathbf{u}^T Z \mathbf{v} \leq 1 \quad (2.10)$$

Solving this last problem is referred to as solving the *polar problem* associated to the function Ω_{θ_γ} .

Before we discuss this auxiliary optimization problem in more detail we first cite an algorithm proposed in [18] and built upon the above results for solving the structured matrix factorization problem. We present the algorithm in the case of the sparse dictionary learning problem.

2.2 Structured matrix factorization algorithm for SDL

We will now present the high-level structure of the SDL matrix factorization algorithm as well as the key elements of each sub-routine, especially the different local optimization strategies.

2.2.1 Meta-algorithm

An adaptation of the proposed structured matrix factorization Meta-Algorithm (Algorithm 1, [18]) for the sparse dictionary learning problem has been implemented by Connor Lane² in MATLAB. In what follows we will present the details of this adapted algorithm and contributed modifications to this implementation.

Algorithm 1: SDL Matrix Factorization

Input: data X , initial factorization (U_{init}, V_{init})
of size r_{init} . **Result:** optimal factorization (U_{final}, V_{final}) of size r_{final}
. **while not converged do**
 1) Local descent to a first-order optimal point (\tilde{U}, \tilde{V}) .
 2) Solve the SDL polar problem at $Z = -\frac{1}{\lambda} \nabla_{\Phi_r} \ell(X, \Phi_r(U, V))$.
 if polar value ≤ 1 then
 | Algorithm has converged.
 else
 | Append the polar solution (u^*, v^*) to (\tilde{U}, \tilde{V}) :
 | $(U, V) \leftarrow ([\tilde{U} \sqrt{\alpha} u^*], [\tilde{V} \sqrt{\alpha} v^*])$ for some $\alpha > 0$.
 end
end

The intuition behind the meta-algorithm in [18] is built upon Corollary 1. Specifically, we begin with an initial solution (U, V) of a certain size r . If such a solution is not a local minimum, then we can perform local descent to arrive at a critical point, which needs to satisfy condition 1 of the corollary. We can then verify if condition 2 is satisfied, in which case we know that we cannot further reduce the objective by increasing r . Alternatively, if condition 2 is not satisfied, then we know that there exists a descent direction (u, v) such that if we augment U by $\sqrt{\alpha}u$ and V by $\sqrt{\alpha}v$, where α is the optimal step size computed in closed form, we can reduce the objective. We can then repeat the process for a factorization of size $r + 1$ until the conditions in Corollary 1 hold, which is guaranteed to happen for finite r , as shown in [18].

²VisionLab/Center for Imaging Science, Johns Hopkins University

2.2.2 Details on local optimization strategies

LOCAL DESCENT

Due to the non-differentiability of the regularization term Θ_γ , we cannot use optimization methods that require smoothness of the cost function such as simple gradient descent. However, it is possible to compute the proximal operator of Θ_γ w.r.t. to U and w.r.t. V , which allows us to define an alternating proximal gradient descent strategy to minimize the following function

$$\begin{aligned} h(U, V) &= \ell(X, \Phi_r(U, V)) + \lambda \Theta_\gamma(U, V) \\ &= \frac{1}{2} \|X - UV^T\|_F^2 + \lambda \sum_{i=1}^r \|U_i\|_2 (\gamma \|V_i\|_1 + (1 - \gamma) \|V_i\|_2) \end{aligned} \quad (2.11)$$

with respect to U and V , i.e. the factorization size r is kept fixed in this step. Using the definition and basic properties of the proximal operator we have:

- For all V , the proximal operator of the continuous convex mapping $U \mapsto \Theta_\gamma(U, V)$ is defined by

$$\begin{aligned} \text{prox}_{\Theta_\gamma(\cdot, V)}(W) &= \arg \min_U \Theta_\gamma(U, V) + \frac{1}{2} \|U - W\|_F^2 \\ &= \sum_i \arg \min_{U_i} \theta_\gamma(U_i, V_i) + \frac{1}{2} \|U_i - W_i\|_2^2 \\ &= \sum_i \text{prox}_{\theta_\gamma(\cdot, V_i)}(W_i), \end{aligned} \quad (2.12)$$

i.e. requiring the computation of the proximal operator of the ℓ_2 norm.

- For all U , the proximal operator of the continuous convex mapping $V \mapsto$

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

$\Theta_\gamma(U, V)$ is defined by

$$\begin{aligned} \mathbf{prox}_{\Theta_\gamma(U, \cdot)}(W) &= \arg \min_V \quad \Theta_\gamma(U, V) + \frac{1}{2} \|V - W\|_F^2 \\ &= \sum_i \arg \min_{V_i} \quad \theta_\gamma(U_i, V_i) + \frac{1}{2} \|V_i - W_i\|_2^2 \\ &= \sum_i \mathbf{prox}_{\theta_\gamma(U_i, \cdot)}(W_i) \end{aligned} \quad (2.13)$$

i.e. requiring the computation of the proximal operator of a weighted sum of ℓ_1 norm and ℓ_2 norm.

From [18, Proposition 5], we know that

$$\mathbf{prox}_{\gamma\|\cdot\|_1 + (1-\gamma)\|\cdot\|_2}(\mathbf{y}) = \mathbf{prox}_{(1-\gamma)\|\cdot\|_2}(\mathbf{prox}_{\gamma\|\cdot\|_1}(\mathbf{y})). \quad (2.14)$$

This relation is useful since we know the proximal operators of the ℓ_1 and ℓ_2 norms:

$$\mathbf{prox}_{\tau\|\cdot\|_2}(\mathbf{y}) = \left(1 - \frac{\tau}{\|\mathbf{y}\|_2}\right)_+ \mathbf{y} \quad (2.15)$$

and

$$(\mathbf{prox}_{\tau\|\cdot\|_1}(\mathbf{y}))_i = \text{sign}(\mathbf{y}_i)(|\mathbf{y}_i| - \tau)_+ \quad \forall i. \quad (2.16)$$

These remarks on the proximal operators $\mathbf{prox}_{\Theta_\gamma(U, \cdot)}$ and $\mathbf{prox}_{\Theta_\gamma(\cdot, V)}$ allow an efficient implementation of an alternating proximal gradient descent algorithm to solve $\min_{U, V} h(U, V)$; at each iteration we are guaranteed a decrease of the objective but this alternation scheme comes without theoretical guarantees on convergence.

APPROXIMATION OF THE POLAR VALUE

Ideally, in the previous step we would use a descent strategy which is guaranteed to converge to a local minimum and not only to a critical point. Once we are at a local minimum we know from Theorem 1 that we only need to check whether this local minimum contains a zero slice (some $i \in [r]$ such that

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

$(U_i, V_i) = (0, 0)$) but unfortunately in the case of non-convex optimization, usual minimization strategies only guarantee convergence to a critical point which can be a local minimum, local maximum or a saddle point. Therefore we need to have an alternative test for optimality which applies to all critical points. From Corollary 1 we know that the polar value tells us whether there exists a descent direction that allows to further decrease the objective or whether we are at the global optimum.

If we substitute in problem (2.7) the definition of the SDL rank-1 regularizer, we get the following problem:

$$\max_{u,v} \quad u^T Z v \quad \text{s.t.} \quad \|u\|_2(\gamma\|v\|_1 + (1-\gamma)\|v\|_2) \leq 1 \quad (2.17)$$

for some fixed matrix Z .

Examining problem (2.17) we remark that we can eliminate one of the variables by noticing that at the optimum we obtain the optimal u^* as $u^* = \frac{Zv^*}{\|Zv^*\|_2}$. Therefore we have the following equivalent problem(s):

$$\max_{v \in \mathbb{R}^N} \quad v^T Z^T Z v \quad \text{subject to} \quad \gamma\|v\|_1 + (1-\gamma)\|v\|_2 \leq 1 \quad (2.18)$$

$$\max_{v \in \mathbb{R}^N} \quad \|Zv\|_2 \quad \text{subject to} \quad \gamma\|v\|_1 + (1-\gamma)\|v\|_2 \leq 1 \quad (2.19)$$

We note that in some cases it reduces to well known problems which we can solve in closed form:

- $\gamma = 0$: $v^* = \text{top eigenvector of } H = Z^T Z$.
- $\gamma = 1$: $v^* = e_{i^*}$ where $i^* = \arg \max_i \|Z_i\|_2$ and e_j denotes an element of the canonical basis of \mathbb{R}^N .

For other cases of γ , i.e. for all $\gamma \in (0, 1)$, however we do not have a closed form solution and need to define an optimization scheme to solve the non-convex polar problem. There are several iterative approximation methods we can define for one of the equivalent formulations above:

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

- *projected power iteration*: given a diagonalizable matrix H and an initial vector v_0 , the power iteration method is defined as:

$$v_{k+1} = \frac{Hv_k}{\|Hv_k\|_2}$$

and converges to the top eigenvector/eigenvalue of H (under some mild assumption on the initial vector v_0). Similarly to [8], we replace the renormalization in the recurrence relation (the projection onto the Euclidean ball) by a projection on our constraint set to obtain the following relation:

$$v_{k+1} = \frac{Hv_k}{\mathcal{P}_{\mathcal{C}}(Hv_k)} \quad (2.20)$$

where $\mathcal{C} = \{v : \gamma\|v\|_1 + (1 - \gamma)\|v\|_2 \leq 1\}$ and the projection onto \mathcal{C} , $\mathcal{P}_{\mathcal{C}}$, is obtained by iteratively solving the proximal operator of the mapping $v \mapsto \gamma\|v\|_1 + (1 - \gamma)\|v\|_2$.

- *(accelerated) proximal gradient ascent*: We can rewrite the polar problem using the indicator function ι :

$$\max_{v \in \mathbb{R}^N} \|Zv\|_2^2 - \iota_{\{\gamma\|v\|_1 + (1-\gamma)\|v\|_2 \leq 1\}} \quad (2.21)$$

and use the proximal gradient algorithm where the proximal operator of $\iota_{\mathcal{C}}$ is simply the projection onto the set \mathcal{C} . An acceleration can be achieved by adding an extrapolation step which translates as described in [29].

- *alternating maximization (or generalized power iteration)*: Instead of deriving the optimal u^* from the result v^* of the optimization over v , we optimize iteratively over both variables and alternate between the two steps

1. $u_{k+1} = \frac{Zv_k}{\|Zv_k\|_2}$
2. $v_{k+1} = \arg \max_v \langle Z^T u_{k+1}, v \rangle$ subject to $\gamma\|v\|_1 + (1 - \gamma)\|v\|_2 \leq 1$,

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

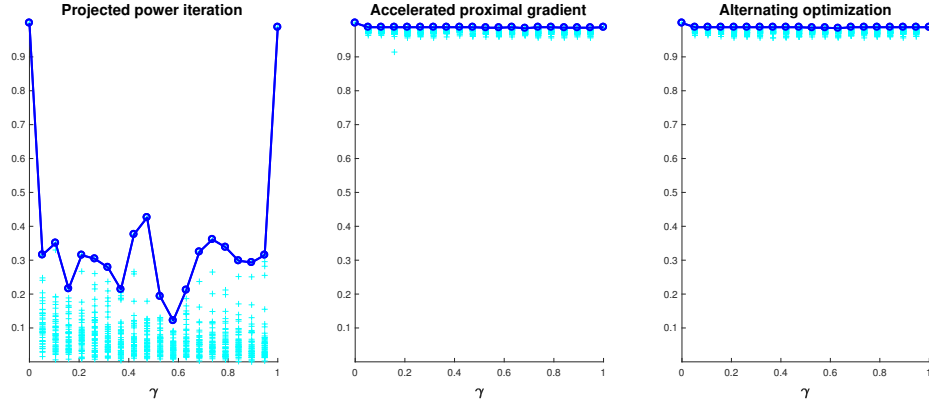
where step 2 is a convex problem which can be solved efficiently by bisection search. This alternating maximization is an adaptation from the method proposed in [21] in the context of sparse principal component analysis.

We compared these three strategies on two types of randomly generated matrices: (1) Gaussian matrix and (2) random matrix with a sparse eigen-vector (different from the top eigen-vector). The results are reported in Figures 2.1 and 2.2 and suggest that the proximal gradient algorithm performs best in both scenarios, even for a small number of random initializations. However, repeating this experiment several times, the alternating maximization yields similar performance as the proximal gradient method. The projected power iteration breaks down for in the case of a Gaussian matrix with dense eigenvectors.

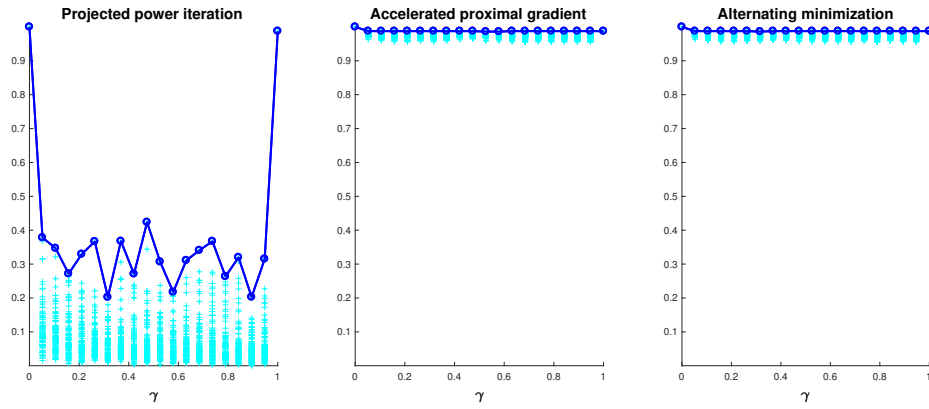
Instead of generating a large number of random initializations, another strategy is to use spectral initialization, i.e. take the eigenvectors of $H = Z^T Z$, since we are trying to solve a generalized eigenvalue problem. We do not claim that the solutions to the generalized eigenvalue problem can be expressed in terms of these eigenvectors, but the empirical results we obtain with this initialization strategy on small-scale problems are satisfactory. Indeed on Figure 2.3 we can see that for both proximal gradient and alternating maximization the spectral initialization leads to similar or better results than using as many random initializations.

Following these results we choose to keep the accelerated proximal gradient and alternating maximization and spectral initialization for our meta-algorithm. Once the polar problem is solved, the polar value is compared to 1: if it is larger, then we append the solution (u^*, v^*) , scaled by a scalar $\alpha^* = \frac{(u^*)^T (X - \tilde{U} \tilde{V}^T) v^* - \lambda \theta_\gamma(u^*, v^*)}{\|u^* (v^*)^T\|_F^2} > 0$ which guarantees a decrease of the objective, to the current factorization.

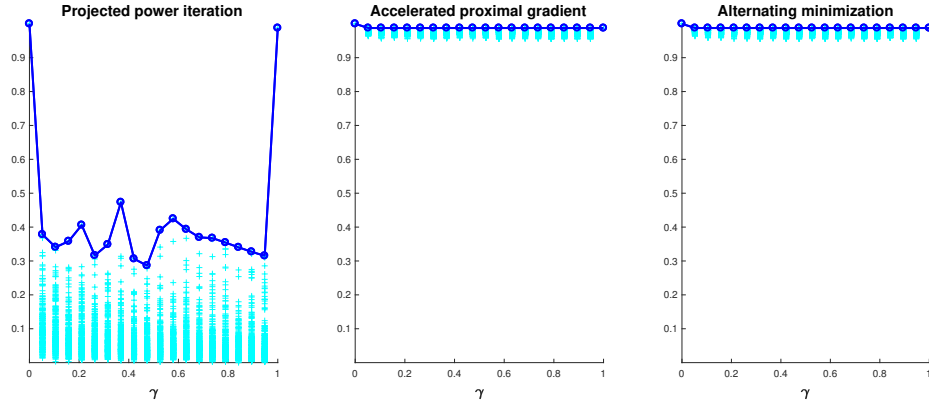
CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM



(a) 50 random initializations.



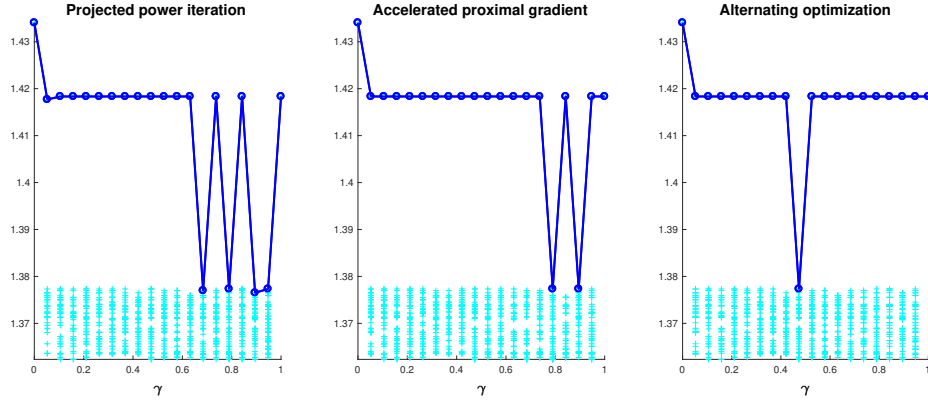
(b) 100 random initializations.



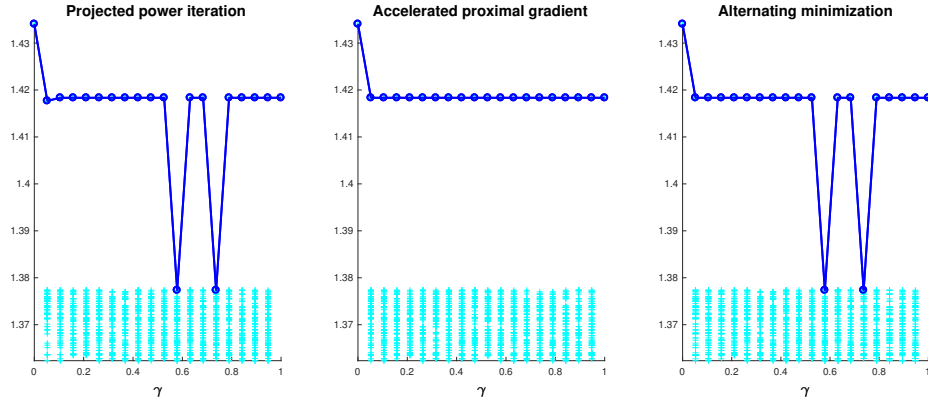
(c) 250 random initializations.

Figure 2.1: Random Gaussian matrix $H = Z^T Z \in \mathbb{R}^{50 \times 50}$ with small condition number (≈ 1.1). The x -axis corresponds to the sparsity trade-off parameter γ . The y -axis gives the value of the cost function $u^T Z v$. All convergence results are represented as cyan crosses, the largest solution is marked by blue circles. Left: projected power iteration. Middle: accelerated proximal gradient. Right: alternating maximization.

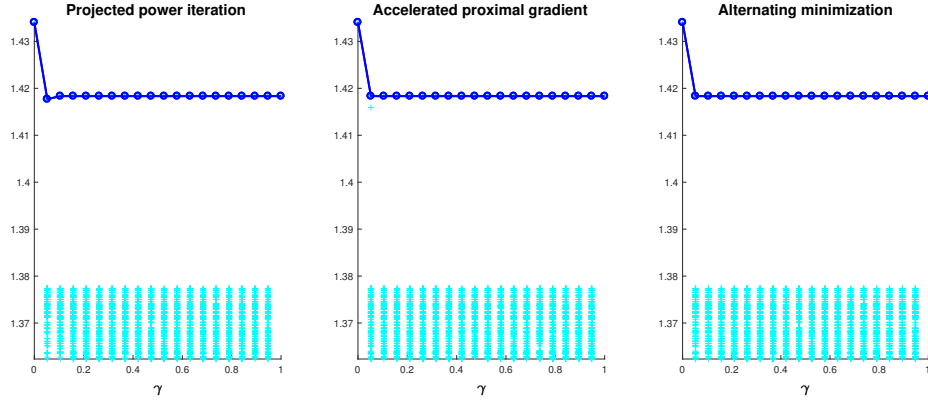
CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM



(a) 50 random initializations.



(b) 100 random initializations.



(c) 250 random initializations.

Figure 2.2: Sparse eigenvector matrix $H = Z^T Z \in \mathbb{R}^{50 \times 50}$ with small condition number (≈ 1.1). The x -axis corresponds to the sparsity trade-off parameter γ . The y -axis gives the value of the cost function $u^T Z v$. All convergence results are represented as cyan crosses, the largest solution is marked by blue circles. Left: projected power iteration. Middle: accelerated proximal gradient. Right: alternating maximization.

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

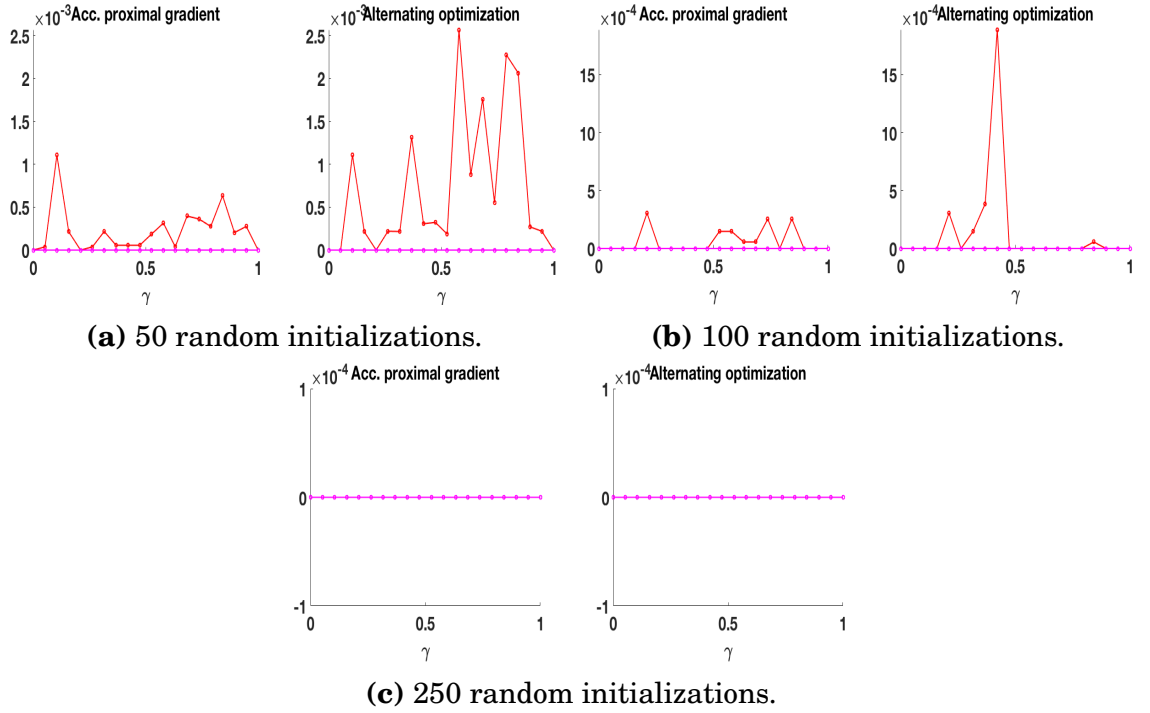


Figure 2.3: Relative difference between the best solution found with spectral initialization (50 eigenvectors of random $H = Z^T Z \in \mathbb{R}^{50 \times 50}$) and the best solution found with random initializations. For each pair of plots, left: accelerated proximal gradient, right: alternating maximization. (red: mean difference; magenta: median difference, taken over 10 experiments).

2.3 Analysis of the polar problem for SDL

In this part we will provide some intuition about the meaning of the polar problem for our overall SDL problem and analyze it as a non-convex optimization problem with a conjectured small number of local minima.

The polar problem that arises in the analysis of the non-convex optimization problem of structured matrix factorization is itself an optimization problem which can be more or less easily solved. Indeed the problem difficulty depends on the rank-1 regularizer θ . For instance, if $\theta(\mathbf{u}, \mathbf{v}) = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2$, then the problem reduces to finding the largest eigenvector of $Z^T Z$ whereas it has been shown that for $\theta(\mathbf{u}, \mathbf{v}) = \|\mathbf{u}\|_\infty \|\mathbf{v}\|_\infty$ the problem is NP-hard [20].

The condition on the polar value can be understood as a general higher-order optimality condition in the sense that if we are at a saddle point (of the optimization problem over (U, V, r)) it allows to find a descent direction; it is a higher order non-smooth saddle point problem. For instance in the low-rank matrix factorization case ($\theta(\mathbf{u}, \mathbf{v}) = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2$) we can show the following result:

Proposition 1. *In the low-rank matrix factorization case,*

$$\min_{U, V, r} f(U, V) = \|X - UV^T\|_F^2 + \lambda \sum_{i=1}^r \|U_i\|_2 \|V_i\|_2 \quad (2.22)$$

the condition $\Omega^\circ(Z) \leq 1$ on the polar function at $Z = -\frac{1}{\lambda} \nabla_{\Phi_r} \ell(X, \Phi_r(U, V))$ is equivalent to the second order optimality condition of the initial problem (2.22).

Proof. Let $r = \text{rank}(X)$. We know that problem (2.22) is equivalent to

$$\min_{U, V, r} \tilde{f}(U, V) = \|X - UV^T\|_F^2 + \frac{\lambda}{2} \sum_{i=1}^r \|U_i\|_2^2 + \|V_i\|_2^2 \quad (2.23)$$

A formal derivation for this equivalence can be found in [5, 7].

We define $W := \begin{bmatrix} U \\ V \end{bmatrix}$ and $h : \mathbb{R}^{(D+N) \times (r+1)} \rightarrow \mathbb{R}_+$ such that $h(W) = \tilde{f}(U, V)$

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

and $\Delta := \begin{bmatrix} \Delta_U \\ \Delta_V \end{bmatrix} \in \mathbb{R}^{(D+N) \times (r+1)}$.

Using standard matrix derivations, we obtain an expression of the Hessian of h : $\nabla^2 h(W)[\Delta, \Delta] = \|\Delta_U V^T + U \Delta_V^T\|_F^2 + 2\langle UV^T - X, \Delta_U \Delta_V^T \rangle + \lambda \|\Delta_U\|_F^2 + \|\Delta_V\|_F^2$.

Let $(\tilde{U}, \tilde{V}) \equiv \tilde{W} \in \mathbb{R}^{(D+N) \times (r+1)}$ be a local minimum of h . The second-order optimality condition is then: $\nabla^2 h(\tilde{W})[\Delta, \Delta] \geq 0$. We know from theorem 1 that it is a global optimum if for some $i \in [r]$, $(\tilde{U}_i, \tilde{V}_i) = (0, 0)$. So the second-order optimality condition reduces to

$$2\langle \tilde{U}\tilde{V}^T - X, \mathbf{u}\mathbf{v}^T \rangle + \lambda \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 \geq 0$$

which is equivalent to the condition on the polar function corresponding to problem (2.23). \square

This result in the Low-Rank Matrix Factorization case, which is a boundary case of our Sparse Dictionary Learning problem (2.5), gives us an intuition of the meaning of the polar problem for the more general case. This general formulation with $\gamma \in (0, 1)$ has no closed-form solution but we found some empirical evidence that the polar problem for this general SDL problem still has a structured optimization geometry as detailed later on.

We recall that we previously derived an equivalent formulation, (2.19), of the initial polar problem (2.17). A first attempt to solve this constrained extrema problem (2.19) approximately is to use one of the previously introduces initialization and optimization strategies and to pick the largest solution. The question however is how difficult is it to reach the global solution, or how many local maxima does this problem have? If this number does not explode exponentially in the dimension one might hope to find (a good approximation of) the global optimum by finding sufficiently many local maxima by local ascent strategies and picking the largest one.

2.3.1 Empirical estimation of the optimization geometry

To get a first intuition of the difficulty of this problem, we study the problem for different values of $\gamma \in (0, 1)$ and on randomly generated matrices Z of rather small dimensions ($D, N \leq 50$). A first observation is that, as expected from the non-convexity of the problem, there always are multiple critical points that the optimization method converges to (see Figures 2.1 and 2.2). To get a better estimation of the number of local optima we run the approximate solvers on $\mathcal{O}(N^2)$ random initializations³. The optimization strategy used to obtain the results in Figure 2.4 is an accelerated proximal gradient ascent, where the proximal operator is simply the projection onto the constraint set $\{v \in \mathbb{R}^N : \gamma \|v\|_1 + (1 - \gamma) \|v\|_2 \leq 1\}$. Alternatively we also used the alternating maximization scheme over u and v which lead to very similar results. We drew the matrix Z in two different ways: (1) from a centered Gaussian distribution with covariance $\Sigma = I$ and (2) such that $Z^T Z$ has a sparse eigenvector that does not correspond to the largest eigenvalue. We also controlled the condition number of the matrix by choosing different values for the reciprocal condition number $rc \in (0, 1]$ (defined by the ratio of the smallest singular value over the largest singular value of the matrix Z). In fact, for the boundary case $\gamma = 0$, we know that the convergence speed of iterative optimization methods, for instance by using the power method, depend on the eigen-gap of $Z^T Z$ and therefore the conditioning of the matrix plays a role. What we observe on Figure 2.4 is that the number of local optima that we find after many random initializations ($2N^2$) appears to be bounded linearly in the dimension N (we removed the sign ambiguity in the reported results, including the sign in the counting leads to an empirical bound of $2N$). We know that for $\gamma = 0$ we have exactly 1 local minimum (the top eigenvector of $Z^T Z$ which is therefore globally min-

³Assuming that there are of the order N local minima, that they are uniformly distributed over the set of feasible points and that they have basins of attraction of similar size, we need $\mathcal{O}(N \log(N))$ initializations to visit all of them at least once, according to the coupon-collector's problem.

imal) and for $\gamma = 1$ there are N local minima (every element of the canonical basis of \mathbb{R}^N). Letting tend γ to each of these boundary cases we observe that the number of detected local minima is consistent with these known results which might suggest that the results for the remaining cases are good estimates of the actual number of local minima. However we cannot exclude the possibility of further local minima which have very small basins of attraction and are therefore unlikely to be reached from an arbitrary point. Additionally the choice of the optimization method might bias the empirical estimation of the optimization geometry. The latter has been pointed out by [16] in the case of under-determined linear regression and separable linear classification problems. A more thorough empirical and theoretic study of the optimization geometry of the studied polar problem is left for future work.

2.3.2 Restricted necessary optimality conditions

With these first empirical results in hand we now go on to the theoretical analysis: first we will give a geometric interpretation of the problem and its solutions, then we will derive an equivalent formulation of the SDL polar problem which on the one hand allows a simpler derivation of second-order optimality conditions and on the other hand allows to draw connections to the well-known problem of sparse PCA.

The objective function of (2.18) is a quadratic function and the constraint is an interpolation of the ℓ_1 and ℓ_2 norm balls as illustrated (in blue) in Figure 2.5. We know that the solutions of this constrained extrema problem lie on the boundary of the constraint. Furthermore, at first-order (FO) optimal points, the constraint boundary is tangent to the level surface. However, since we want to estimate the number of local maxima and the problem being non-convex, we are interested in second-order (SO) optimal points. The SO optimality condition depends on the curvature of the surface of level and constraint set and informally it requires the constraint surface to have strictly larger curvature at a first-order optimal point than the level surface for it to be second-order

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

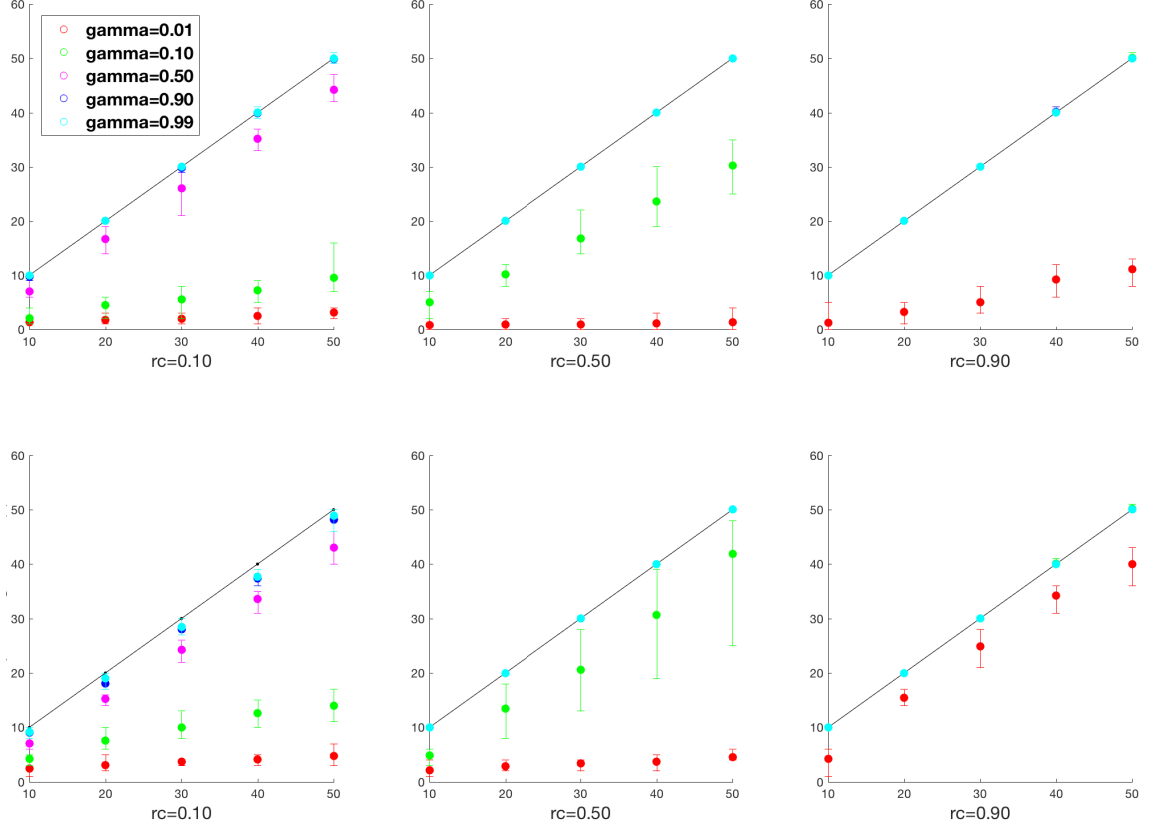


Figure 2.4: Estimated number of local maxima of the SDL polar problem (y -axis) for different problem sizes as a function of the dimension N (x -axis). Top: $H = Z^T Z$ where each $Z_i \sim \mathcal{N}(0, I)$. Bottom: H is generated such that it has at least one sparse eigenvector using a similar strategy as [21]. For each plot, 10 matrices with a fixed condition number were generated per dimension and the problem solved by accelerated proximal gradient ascent. Left to right column: results for matrices with reciprocal condition number equal to 0.1, 0.5 and 0.9 respectively.

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

optimal. An example in \mathbb{R}^2 is given in Figure 2.5. Due to the non-differential ℓ_1 term in the constraint and the boundary case $\gamma = 1$ we have to take some precautions to formalize this informal characterization which uses the surfaces' curvatures. Before we give the formal expression of the FO and SO optimality conditions we derive an equivalent formulation of the SDL polar problem in Proposition 2.

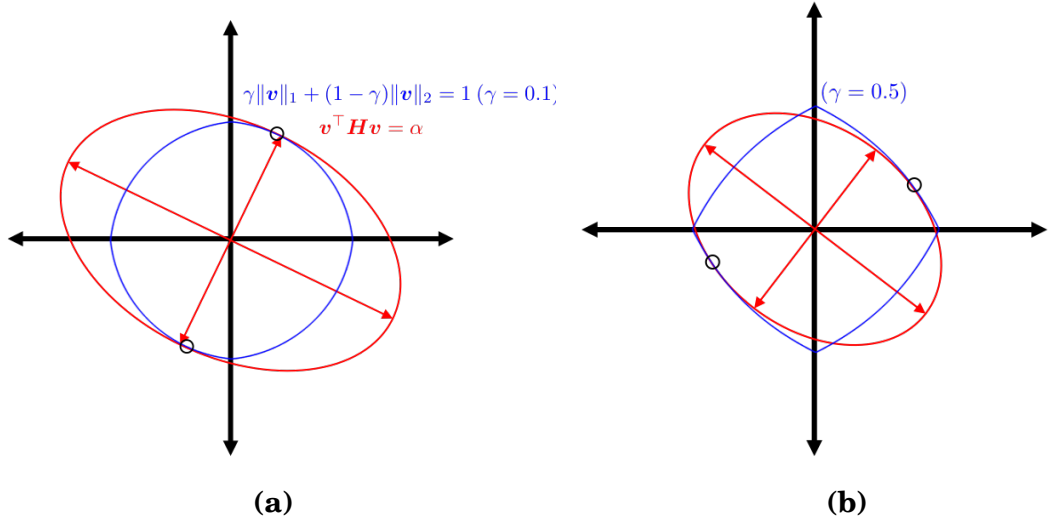


Figure 2.5: Level and constraint sets of an SDL polar problem in \mathbb{R}^2 . FO optimal points are indicated by small circles. (a) The FO optimal points are also second-order optimal; (b) The FO optimal points violate the SO optimality condition.

Proposition 2. *The sparse dictionary learning polar problem is equivalent to solving the optimization problem*

$$\min_{\substack{\mathbf{v} \in \mathbb{R}^N \\ \|\mathbf{v}\|_2=1}} \frac{\gamma\|\mathbf{v}\|_1 + 1 - \gamma}{\|Z\mathbf{v}\|_2} \quad (2.24)$$

where $Z = -\frac{1}{\lambda} \nabla_{\Phi} \ell(X, \Phi(U, V)) \in \mathbb{R}^{D \times N}$.

Proof. We start with formulation (2.19) of the polar problem derived earlier. Using a result on constrained extrema that states that the solutions of constrained maximization of a convex function are located at the boundaries of

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

the constraint set we have:

$$\max_{\mathbf{v} \in \mathbb{R}^N} \frac{\|Z\mathbf{v}\|_2}{\gamma\|\mathbf{v}\|_1 + (1-\gamma)\|\mathbf{v}\|_2} \quad \text{subject to } \gamma\|\mathbf{v}\|_1 + (1-\gamma)\|\mathbf{v}\|_2 = 1 \quad (2.25)$$

Assuming that $\text{null}(Z) \neq \mathbb{R}^N$, we can add the constraint $\mathbf{v} \notin \text{null}(Z)$. The positivity of the objective function and invariance to the scaling of \mathbf{v} allow us to write

$$\min_{\mathbf{v} \in \mathbb{R}^N} \frac{\gamma\|\mathbf{v}\|_1 + (1-\gamma)\|\mathbf{v}\|_2}{\|Z\mathbf{v}\|_2} \quad \text{subject to } \|\mathbf{v}\|_2 = 1 \quad (2.26)$$

□

The objective function of (2.26) is a convex combination of the objective functions of the variational definition of the $\|\cdot\|_{1,2}$ and $\|\cdot\|_{2,2}$ operator norms on Z . We point out that [19] use this formulation to compute approximate solutions for the well-known sparse PCA problem [39] by using an algorithm designed to solve non-linear eigenvalue problems (the problem being non-linear in the eigenvector). This connection to the sparse PCA problem provides an additional motivation to study the SDL polar problem more deeply to understand its optimization landscape.

We will now present the derivations of first- and second-order optimality conditions for problem (2.24) to further develop the previously mentioned geometric intuitions for this polar problem. In order to eliminate problems due to the non-differentiable ℓ_1 term, we compute the gradient and Hessian of the objective function for a fixed support $\sigma \subseteq [N]$ and signature $\mathbf{s} \in \{-1, 0, 1\}^N$. Let $Z \in \mathbb{R}^{D \times N}$ and $\mathbf{v} \in \mathbb{R}^N$ with support σ and signature \mathbf{s} , without loss of generality we can assume $\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{v}}_1 \odot \mathbf{s}_1 \\ \mathbf{0} \end{bmatrix}$ where $\mathbf{s} = \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{0} \end{bmatrix}$. We also partition matrix Z accordingly: $Z = [Z_1 \ Z_2]$. We denote $S_1 = \text{diag}(\mathbf{s}_1)$, $H = Z^T Z = \begin{bmatrix} H_{11} & H_{12} \\ H_{21} & H_{22} \end{bmatrix}$ and $\tilde{H}_{11} = S_1 Z_1^T Z_1 S_1$. With these notations in hand, we give the following expression for the first-order and second-order optimality conditions:

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

Proposition 3. *Let $Z \in \mathbb{R}^{D \times N}$ and $v \in \mathbb{R}^N$ with support σ and signature s . v is a local minimum of (2.24) if it satisfies the following conditions:*

$$1) \quad (FO_\sigma) \quad \tilde{H}_{11}\tilde{v}_1 = \frac{\eta}{\kappa}(\gamma\mathbf{1} + (1-\gamma)\tilde{v}_1) \quad (2.27)$$

$$2) \quad (FO_{\sigma^c}) \quad H_{21}S_1\tilde{v}_1 \in \left[-\frac{\eta\gamma}{\kappa}, \frac{\eta\gamma}{\kappa}\right]^{N-|\sigma|} \quad (2.28)$$

$$3) \quad (SO_\sigma) \quad \eta \left(\frac{\gamma^2}{\kappa^2}(\mathbf{1}^T \xi)^2 + \frac{1-\gamma}{\kappa} \right) \geq \xi^T \tilde{H}_{11} \xi \quad (2.29)$$

for all $\xi \in \mathcal{T}_{\tilde{v}_1} \mathbb{S}^{|\sigma|-1}$ such that $\|\xi\|_2 = 1$

where we defined $\eta = \tilde{v}_1^T \tilde{H}_{11} \tilde{v}_1$ and $\kappa = \gamma(\mathbf{1}^T \tilde{v}_1) + (1-\gamma)$. Given a manifold \mathcal{M} , $\mathcal{T}_u \mathcal{M}$ denotes the tangent space at $u \in \mathcal{M}$ to the manifold.

The proof of this proposition uses notions and results of differential geometry and manifold optimization from [1] and is given in the Appendix.

The FO necessary condition (FO_σ) expresses the previous geometric observation that at a first-order optimal point, the normals to the surface of the constraint and level sets are parallel. And in the boundary case $\gamma = 0$ we recover the result that the solution to (2.24) is given by the top eigen-vector of H . However for other choices of γ , a geometric interpretation of condition (SO_σ) in terms of curvature of the surface of level and constraint sets is less straightforward. But these optimality conditions on the (Riemannian) gradient and Hessian could possibly be used to bound the number of local optima of the SDL polar problem, by invoking a similar proof as in [14] for the problem of tensor decomposition. This direction is left for future work.

2.4 Extension to discriminative sparse dictionary learning

In this section we present a potential extension of our SDL formulation to discriminative SDL and explain how the theoretical results and our algorithm

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

easily extend to this new problem, using the fact that the matrix factorization mapping Φ_r can be chosen different from the regular matrix-matrix product (UV^T).

Learning a sparse dictionary and decomposition directly from the data has several different applications, for instance classification. In order to classify the data, a possibility is to consider the sparse codes as a nonlinear transformation of the data and to learn a classifier on these features. Indeed it is admitted that sparse codes are well-suited for classification of static data ([24] and [33] also show how they can be used for classification of time-series data). Since these features do not depend on some pre-defined dictionary but on a data-dependent dictionary it seems natural to “design” the dictionary and the sparse codes for the classification task, i.e. to make them discriminative. Given a training set with known class labels, this can be achieved by adding a discriminative loss term to the model, as it has been proposed in [26]. The authors call their model *Supervised Dictionary Learning* to emphasize the explicit use of the class labels during dictionary learning as opposed to an unsupervised dictionary learning approach.

In order to integrate this idea of learning discriminative sparse codes into the SDL matrix factorization framework we assume that the data can be separated by a linear classifier in the latent sparse coding coefficient space. Consider that there are L different classes and that for each data point \mathbf{x}_i we know its class label y_i ; without loss of generality we can assume $y_i \in [L]$ for all $i \in [N]$. Given the sparse code $V_{<i>} \in \mathbb{R}^r$ of \mathbf{x}_i and a classifier $\{W, \mathbf{b}\} \in \mathbb{R}^{L \times r} \times \mathbb{R}^L$, the predicted label \hat{y}_i can be defined via the *softmax* function (assuming $\hat{y}_i = \arg \max_{l \in [L]} W_{<l>}^T V_{<i>} + \mathbf{b}_l$). Under these assumptions we can naturally use the cross-entropy loss on the softmax classifier in our objective function:

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

$$\begin{aligned}
\min_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r} \\ W \in \mathbb{R}^{L \times r}, r \in \mathbb{N} \\ b \in \mathbb{R}^L}} & \underbrace{\frac{1}{2} \|X - UV^T\|_F^2 + \kappa \sum_{i=1}^N \mathcal{K}(y_i, f(V_{<i>}, \{W, \mathbf{b}\}))}_{=\ell((X, y), \Phi_r(U, V, W), \mathbf{b})} \\
& + \underbrace{\lambda \sum_{j=1}^r \left\| \begin{bmatrix} U_j \\ W_j \end{bmatrix} \right\|_2 (\gamma \|V_j\|_1 + (1 - \gamma) \|V_j\|_2)}_{=\Theta_\gamma(U, V, W)}
\end{aligned} \tag{2.30}$$

where $f(\alpha, \{W, \mathbf{b}\}) = W_{<l>}^T \alpha + \mathbf{b}_l$ is the score function and \mathcal{K} the cross-entropy loss.

This model fits into the matrix factorization framework with the following elements:

- The elemental mapping

$$\begin{aligned}
\phi : \mathbb{R}^D \times \mathbb{R}^N \times \mathbb{R}^L &\rightarrow \mathbb{R}^{(D+L) \times N} \\
(\mathbf{u}, \mathbf{v}, \mathbf{w}) &\mapsto \begin{bmatrix} \mathbf{u}\mathbf{v}^T \\ \mathbf{w}\mathbf{v}^T \end{bmatrix}
\end{aligned} \tag{2.31}$$

- The r -element factorization mapping

$$\begin{aligned}
\Phi_r : \mathbb{R}^{D \times r} \times \mathbb{R}^{N \times r} \times \mathbb{R}^{L \times r} &\rightarrow \mathbb{R}^{(D+L) \times N} \\
(U, V, W) &\mapsto \sum_{i=1}^r \phi(U_i, V_i, W_i)
\end{aligned} \tag{2.32}$$

- The rank-1 regularizer

$$\begin{aligned}
\theta_\gamma : \mathbb{R}^D \times \mathbb{R}^N \times \mathbb{R}^L &\rightarrow \mathbb{R}_+ \\
(\mathbf{u}, \mathbf{v}, \mathbf{w}) &\mapsto \left\| \begin{bmatrix} \mathbf{u} \\ \mathbf{w} \end{bmatrix} \right\|_2 (\gamma \|\mathbf{v}\|_1 + (1 - \gamma) \|\mathbf{v}\|_2)
\end{aligned} \tag{2.33}$$

- The loss function $\ell((X, y), Z, \mathbf{b})$ is differentiable and jointly convex w.r.t.

CHAPTER 2. GLOBAL OPTIMALITY FOR THE SDL PROBLEM

(Z, b) . (Here we consider that the bias b as an auxiliary variable which is not part of the factorization; the theoretical framework of our main reference [17] covers this case, but for simplicity we omitted this detail in theorem 1).

Proposition 4. *Solving the polar problem for the sparse dictionary learning problem (2.5) is equivalent to solving the polar problem for the discriminative sparse dictionary learning (2.30).*

Proof. If we define $\tilde{u} = \begin{bmatrix} u \\ w \end{bmatrix}$ and adjust the problem dimensions, then we obtain the same formulation of the SDL polar problem as in (2.17). \square

A motivation for integrating the classifier in the sparse dictionary learning matrix factorization framework is twofold: this additional variable does not make the optimization scheme more difficult to solve than the regular SDL formulation as shown above and it allows to make statements about the optimal size for the dictionary trained for a discriminative (and also reconstructive) task. Indeed, when reviewing the literature on classification via sparse coding the heuristics for the dictionary size vary from under-completeness (to prevent over-fitting, for instance [26]) to larger dictionaries (if the dictionary is also required to be reconstructive [31]) and the dictionary size is then also chosen via cross-validation which can come at a high computational cost.

However, empirically validating our model on synthetic data and developing some intuition on the formulation's sensitivity to changing parameters turned out to be challenging because of the difficulty to generate data such that we know the joint ground truth $\{U^*, V^*, W^*, b^*\}$, or at least a good proxy of it. Furthermore the impact of the inexact polar solutions (except for the boundary cases $\gamma \in \{0, 1\}$) is not clear yet and still needs to be understood in the regular SDL case. And we conjecture that, in order to obtain interpretable results, we need to add a positivity constraint on the coding variables to guarantee a linear classifier W on these variables.

3 Subspace Clustering and Sparse Dictionary Learning

In this chapter we review two state-of-the-art methods for subspace clustering, belonging to the broader class of spectral methods. They promote different structures to provide compact representations of the data (sparsity or low-rank) which allow to recover a representation of the low-dimensional structure of the data drawn from multiple subspaces by finding a specific decomposition over a self-expressive dictionary (i.e. the data matrix serves as dictionary). As pointed out in the introduction, both sparse dictionary learning and subspace clustering assume the same data model and we will now position our SDL model relatively to the two previously mentioned subspace clustering methods and demonstrate how subspace clustering can be achieved with our formulation while providing additional information on the optimality of the dictionary size w.r.t. the generating factorization size.

3.1 Subspace clustering by sparse or low-rank representation

Given a large amount of high-dimensional data, stacked as columns of a matrix $X \in \mathbb{R}^{D \times N}$ (for instance D could be the number of pixels in an image and N the number of images), it is often preferable to invoke a low-dimensional representation of the data or to cluster the data such that each cluster is well approximated by a low-dimensional subspace in order to perform different data processing and learning tasks more efficiently. In computer vision for instance these could be denoising, classification or inpainting.

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

If one assumes that the data is well represented by a single low-dimensional subspace, then a classical approach to dimensionality reduction is Principal Component Analysis (PCA). This method and its variants perform well as long as the data can indeed be approximated by a single subspace. But if the structure becomes more complex, i.e. the data rather belongs to multiple subspaces as illustrated in Figure 3.1 and is potentially corrupted by noise or outliers, more advanced techniques of fitting the data to such low-dimensional structure are required. Given such a collection of datapoints coming from multiple subspaces, the subspace clustering problem consists of recovering the following information: number and dimensions of the subspaces, a basis for each subspace and the segmentation of the data, ideally without requiring any prior on the number and dimensions of the subspaces.

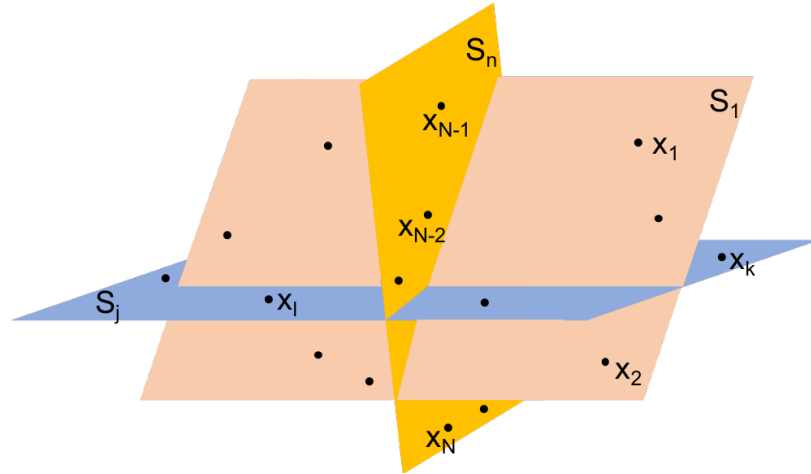


Figure 3.1: Data $X = [x_1 \dots x_N]$ lies on a union of linear low-dimensional subspaces $\mathcal{S} = (S_1 \cup \dots \cup S_n) \subset \mathbb{R}^D$.

One way to tackle this problem is to find a self-expressive representation of the data as illustrated in Figure 3.2 and to use this representation to infer the clustering of the points w.r.t. the different subspaces: given $X \in \mathbb{R}^{D \times N}$, one can assume the self-expressiveness property and seek a matrix $C \in \mathbb{R}^{N \times N}$ such that $X = XC$. To prevent the trivial solution $C = I$ one needs to enforce $\text{diag}(C) = 0$, which does not contribute to the aim of recovering the low-

dimensional structure of the data, one enforces specific constraints onto this matrix C which translate different assumptions on the form of the data, for instance sparsity and low-rank. We will consider these two different situations as proposed in [11] and [12, 23] respectively.

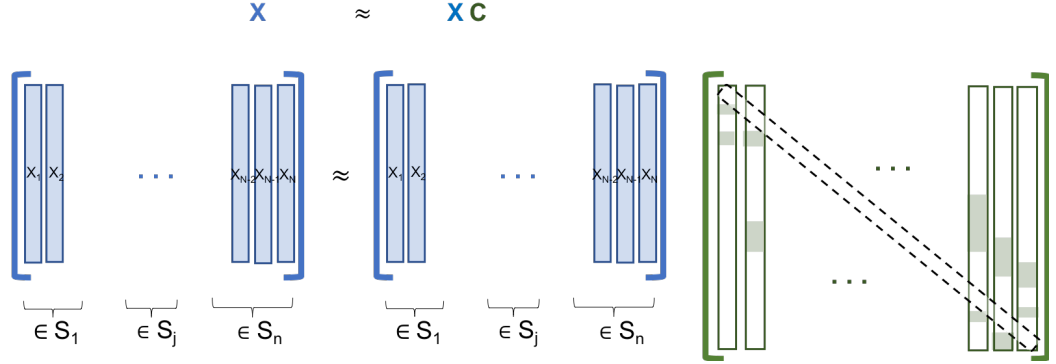


Figure 3.2: Data $X = [x_1 \dots x_N]$ serves as a self-expressive dictionary. The coefficient matrix C is constrained to be zero on its main diagonal.

3.1.1 Sparse subspace clustering

The idea behind the sparse subspace clustering formulation (SSC) is the following: a datapoint x can be approximated by a linear combination of the other datapoints from the set. Such a representation is not unique – there are infinitely many solutions –, but if one requires it to be sparse and if the subspaces are well arranged then it is likely that only points lying in same subspace as x will be selected to represent the given point or at least most of the selected points will belong to the same subspace as x . Such a representation is referred to as *subspace-sparse*. Therefore one gets to solve the following problem:

$$\min \|C\|_0 \quad \text{subject to } X = XC \text{ and } \text{diag}(C) = 0 \quad (3.1)$$

where the second constraint prevents the trivial solution $C = I$, i.e. each point being represented by itself. This problem is shown to be NP-hard due to the pseudo-norm ℓ_0 but the classical tight relaxation using the sparsity inducing ℓ_1

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

norm allows to solve a convex problem [11]:

$$\min \|C\|_1 \quad \text{subject to } X = XC \text{ and } \text{diag}(C) = 0 \quad (3.2)$$

The above formulation requires an exact reconstruction of the data ($X = XC$) but if one assumes that the data is contaminated by some random noise E , then it is preferable to minimize the reconstruction error $\|X + E - XC\|$, which gives the extension of SSC to the presence of noise:

$$\min \|C\|_1 + \frac{\tau}{2} \|E\|_F^2 \quad \text{subject to } X = XC + E \text{ and } \text{diag}(C) = 0 \quad (3.3)$$

The authors of [11] show that under appropriate assumptions on the distribution of the data and the configuration of the subspaces the solution of the above formulation is subspace-sparse and can therefore be used to build a similarity graph to deduce the segmentation of the data.

3.1.2 Low-rank subspace clustering

Instead of arguing in terms of sparse representation, one can consider that if the subspaces are well arranged such that the rank r of the data is equal to the sum of the dimensions of the subspaces, therefore potentially $r < \min\{D, N\}$ (if D and N are indeed sufficiently large). This leads to the following non-convex problem

$$\min \text{rank}(C) \quad \text{subject to } X = XC \quad (3.4)$$

which can be relaxed into a convex problem by noticing that similarly to the ℓ_0/ℓ_1 substitution, the nuclear norm $\|\cdot\|_*$ is a convex envelope for the matrix rank function¹.

If in addition one assumes the data to be corrupted by random noise, i.e. X is the sum of the clean data A and the noise E , the following convex problem

¹Given a square or rectangular matrix M the nuclear norm is defined as the sum of the singular values of M , i.e. $\|M\|_* = \sum_{i=1}^r \sigma_i$, where the singular values are obtained from the singular value decomposition of M : $M = U \text{diag}(\sigma) V^T$.

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

can be solved to get a subspace-sparse representation of the data under the appropriate assumptions.

$$\min \quad \|C\|_* + \frac{\tau}{2}\|E\|_F^2 \quad \text{subject to } A = AC \text{ and } X = A + E \quad (3.5)$$

Under appropriate assumptions the authors of [12] show that for clean data and in the presence of noise the problem can be solved in closed form.

These two methods both yield state-of-the-art results for subspace clustering for clean as well as corrupted data and only require fairly simple tools from optimization and linear algebra. But their major drawback lies in their complexity in the number of signals. And the self-expressiveness assumption, while leading to efficient clustering methods, does not provide any information on the optimality of the decomposition w.r.t. the set of sparse decompositions $\{(D, C, r) \in \mathbb{R}^{D \times r} \times \mathbb{R}^{r \times N} \times \mathbb{N}^* : X = DC\}$ since the self-expressiveness assumption implies the use of the fairly large dictionary $D = X$.

Given our analysis from the previous chapters we can therefore ask whether our SDL matrix factorization framework can provide such information for the problem of subspace clustering. In other words we are interested in the “compactness” of the factorization in the problem setting of subspace modeling and clustering.

3.2 Subspace clustering by sparse dictionary learning

The underlying assumptions for sparse dictionary learning are to some extent related to those of subspace clustering: by attempting to sparsely decompose the data over some (overcomplete) set of basis elements one implicitly assumes that the data lies in a union of linear low-dimensional subspaces which allows to express each signal as a linear combination of some basis elements of

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

the subspace it belongs to, in other words one attempts to recover a subspace-sparse representation. Hence, one choice for the dictionary is to form it as the concatenation of such basis elements leading to subspace-preserving codes (decomposition coefficients) which are used to derive a similarity graph and the segmentation of the data. The data model and its approximated decomposition are illustrated in Figures 3.1 and 3.3.

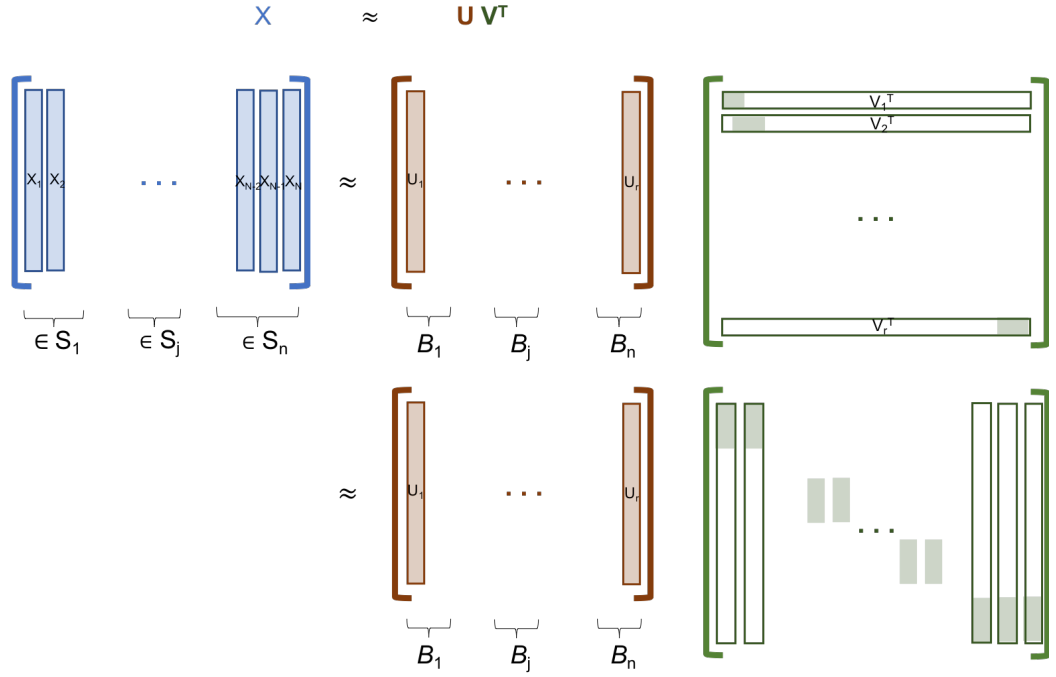


Figure 3.3: Data X expressed as (sparse) linear combinations of subspace specific basis elements ($B_j, j \in [n]$).

The two previously presented models, SSC and LR-SC, can be considered as special cases of this general formulation, where the signals are considered as a concatenation of (overcomplete) subspace bases. The question now is whether one can recover a more compact and still subspace-preserving decomposition of the data by jointly learning the dictionary and decomposition coefficients without fixing the factorization size in advance. In other words, is there a matrix factorization regularization function and associated structured matrix factorization problem which allows to capture the specific low-dimensional structure

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

of the data and to find an optimal generating factorization size? If one chooses to drop the self-expressiveness in (3.3) and (3.5) in order to find an optimal subspace-sparse representation of the data, one can obtain the following formulations²:

$$\min_{\substack{U, V \\ \|U_i\|_2 \leq 1, \forall i}} \|V\|_1 + \frac{\tau}{2} \|X - UV^T\|_F^2 \quad (3.6)$$

$$\min_{\substack{U, V \\ \|U_i\|_2 \leq 1, \forall i}} \|V\|_* + \frac{\tau}{2} \|X - UV^T\|_F^2 \quad (3.7)$$

For problem (3.7), since we want to find a compact factorization of the data into U and V , instead of assuming only V to be low-rank, we can assume UV^T to have low-rank, leading to

$$\min_{U, V, r} \|UV^T\|_* + \frac{\tau}{2} \|X - UV^T\|_F^2 \quad (3.8)$$

We can now establish the link to our SDL formulation by using the well-known variational form of the nuclear norm: given a matrix $M \in \mathbb{R}^{D \times N}$ two variational forms of its nuclear norm are defined as:

$$\|M\|_* = \inf_r \inf_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r} \\ X = UV^T}} \sum_{i=1}^r \|U_i\|_2 \|V_i\|_2 \quad (3.9)$$

$$= \inf_r \inf_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r} \\ X = UV^T}} \frac{1}{2} \sum_{i=1}^r \|U_i\|_2^2 + \|V_i\|_2^2 \quad (3.10)$$

This alternative definition of the nuclear norm allows us to recover one boundary case of our SDL model:

$$\min_{\substack{U \in \mathbb{R}^{D \times r} \\ V \in \mathbb{R}^{N \times r}, r \in \mathbb{N}}} \frac{1}{2} \|X - UV^T\|_F^2 + \lambda \sum_{i=1}^r \|U_i\|_2 \|V_i\|_2 \quad (3.11)$$

²As before, we add the norm constraint on U to prevent unbounded solutions.

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

which is equivalent to the rank minimization problem

$$\min_{Y \in \mathbb{R}^{D \times N}} \frac{1}{2} \|X - Y\|_F^2 + \lambda \|Y\|_*, \quad (3.12)$$

as showed in [30]. This convex low-rank minimization problem is known to have a closed form solution. But it does not provide any information on the data segmentation.

Similarly for problem (3.6), we use again the remark that constraining the norm of the columns of U is equivalent to penalizing their sum:

$$\min_{U, V, r} \frac{1}{2} \|X - UV^T\|_F^2 + \lambda \sum_{i=1}^r \|U_i\|_2 \|V_i\|_1 \quad (3.13)$$

which corresponds to the second extreme case of our problem (2.5) ($\gamma = 1$). Here again the varying factorization size leads to a (trivial) solution ($U = X, V = I$) which does not give any insight into the underlying structure of the data.

But since we assume that there exists a sparse decomposition for the data X (by construction), we can interpolate these two problems to attempt to obtain a compact (low-rank) and subspace-sparse factorization. In other words, since our SDL matrix factorization formulation also involves optimizing over the factorization size, by choosing a suitable regularization function, i.e. a suitable choice for the sparsity parameter γ , one can attempt to recover the number of subspaces, their dimensions and a basis for each subspace in form of a dictionary as well as the representation of the data over this dictionary, i.e. the segmentation of the data collection. If we denote the dimension of subspace S_j by d_j and if we have $\sum_{j=1}^n d_j < N$, then we expect the resulting factorization of the data to be more compact as opposed to the self-expressive dictionary from SSC, LRR or LR-SC.

3.3 Experiments

With this framing of our SDL formulation w.r.t. two existing subspace clustering methods we now proceed to the verification of our claim on the recovery of a compact subspace-preserving factorization of the data, i.e. for some value for γ can we perform subspace clustering while obtaining an optimal factorization size?

First we will explain our data model, discuss our initialization strategy and stopping criteria which are necessary due to the use of approximate polar solutions and finally we compare the results of our approach to existing subspace clustering methods.

3.3.1 Preliminaries

SYNTHETIC DATA

We analyze our method by applying it on synthetic data: We generate signals according to two different models: (1) the collection of subspaces is independent³ and (2) the subspaces are disjoint⁴. Requiring the subspaces to be independent is a strong assumption but intuitively it makes the subspace estimation and clustering easier. We choose these two data models since the methods we are comparing our model to (SSC and LR-SC) give guarantees of subspace recovery for these cases.

For each experiment we generate 100 signals per subspace in the ambient space \mathbb{R}^{30} . Each subspace is of dimension \mathbb{R}^3 . For model (1) we generate signals from 9 independent subspaces, for model (2) we draw 12 subspace bases $\{U^j\}_{j=1}^{12}$ at random from a normal distribution⁵.

³A collection of subspaces $\{S_j\}_{j=1}^n$ is said to be independent if $\dim(\bigoplus_{j=1}^n S_j) = \sum_{j=1}^n \dim(S_j)$ where \bigoplus denotes the direct sum operator. [11]

⁴A collection of subspaces $\{S_j\}_{j=1}^n$ is said to be disjoint if every pair of subspaces intersect only at the origin. [11]

⁵The generated subspaces are disjoint with probability 1.

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

REMARK ON THE CHOICE OF INITIALIZATION AND STOPPING CRITERION

Assuming that every step of our algorithm 1 is solved exactly, the algorithm is guaranteed to converge to a global solution, starting from any initial factorization. But as pointed out in Part 2 we cannot solve the polar problem exactly and we therefore have to take this approximation error into account. A phenomenon we observe throughout our experiments (on synthetic and real data) is that the final factorization produced by the algorithm always contains a lot of “small” factors with a vanishing rank-1 regularizer θ_γ and some factors with much larger θ_γ (see Figure 3.4 for an example of the evolution of the ratio between the smallest and largest rank-1 regularizer throughout the meta-iterations of our algorithm using approximate polar solutions). This could be due to numerical issues in any of the subroutines of the overall algorithm⁶ but it could also be the case that the global solution for our problem (2.5) indeed contains many small rank-1 factors. An argument for the former is that when looking at the low-rank matrix factorization case ($\gamma = 0$) for some matrix X which is exactly low rank, the algorithm recovers the optimal solution (which we also know in closed form), starting from an empty factorization, i.e. it optimizes over the successive polar directions. But if we corrupt each solution of the polar problem with some small white Gaussian noise, we obtain a similar behavior as in Figure 3.4, i.e. a very large factorization containing many small additional factors as reported in Figure 3.5. In parallel we can compare the evolution of the polar and objective functions, see Figure 3.6: in all cases the polar and the objective functions have two stages of evolution: a fast decrease until the factorization size reaches r_{data} and a very slow decrease throughout the rest of the optimization.

These observations suggest that instead of defining the convergence criterion for our algorithm on the polar value (*converged* $\equiv (\Omega_{\theta_\gamma}^\circ(Z) = 1)$) which is likely to lead to an overestimation of the factorization size, we rather need to

⁶In Section 2.2 we empirically convinced ourselves that the non-convex polar problem has a nice optimization landscape and that spectral initialization leads to good estimations of the true solution.

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

identify the “knee” in either the polar function or the objective which occur at the generating factorization size r_{data} .

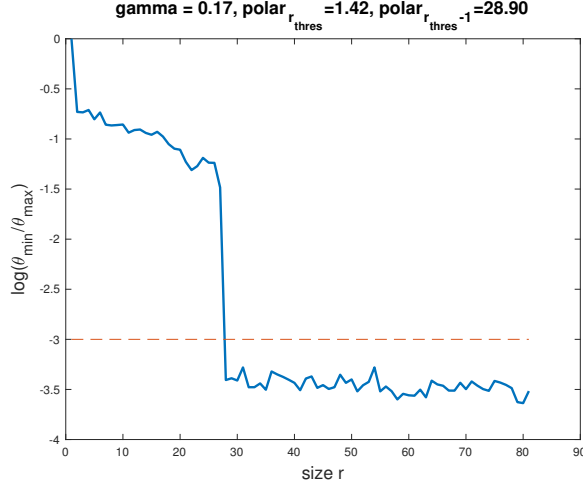


Figure 3.4: Ratio of smallest and largest rank-1 regularizer throughout the meta-iterations of algorithm 1, initialized with the empty factorization. The x -axis represents the factorization size r . We stop the algorithm once the polar value drops below 1.1. Problem dimensions: $X = UV^T \in \mathbb{R}^{30 \times 900}$ and $\text{rank}(X) = 27$, $\gamma = 0.17$. The red dotted line corresponds to $y = \log(0.001)$. The ratios and reference value y are reported in log-scale.

FORWARD SDL MF

Our algorithm 1 is designed such that it can be initialized with any factorization; for instance we could use initialization strategies used in other approaches to solve the problem: initialize the dictionary U with a random subset of the data matrix X and random sparse coefficients V . However we choose to initialize the algorithm with a zero factorization and the algorithm will therefore run for at least r^* iterations, where r^* is the factorization size of the global solution. This specific choice of initialization also appears in a similar algorithmic setting in [4]. We will refer to this choice zero initialization for our algorithm as *forward SDL MF*.

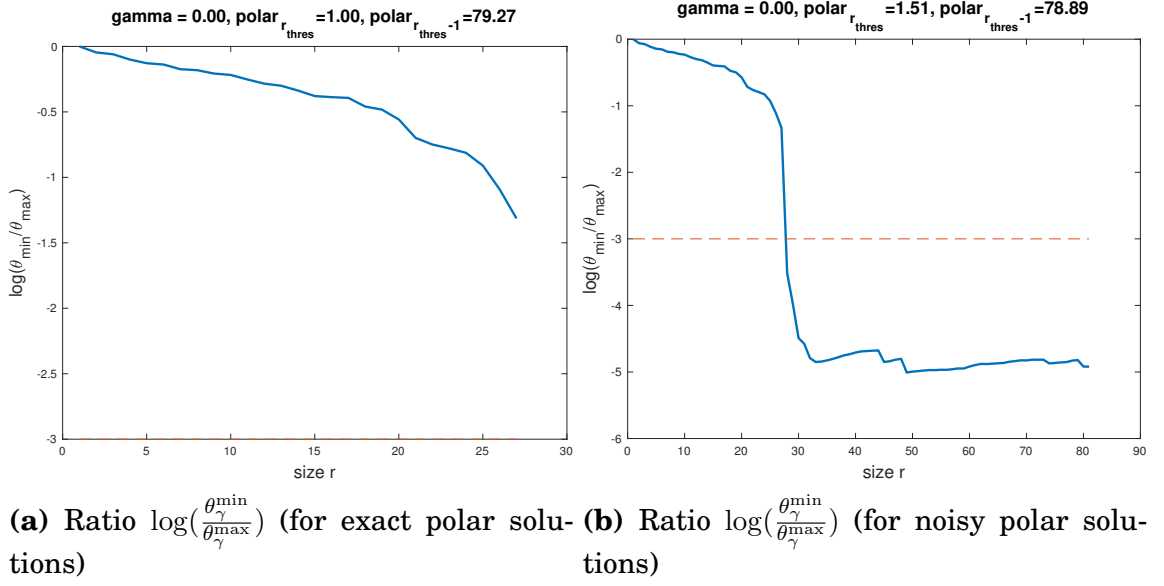


Figure 3.5: Ratio of smallest and largest rank-1 regularizer throughout the meta-iterations of algorithm 1, initialized with the empty factorization. The x -axis represents the factorization size r , which coincides with the meta-iterations in this case. We stop the algorithm once the polar value drops below 1.1. Problem dimensions: $X = UV^T \in \mathbb{R}^{30 \times 900}$ and $\text{rank}(X) = 27$. The red dotted line corresponds to $y = \log(0.001)$.

3.3.2 Recovering the generating factorization size

The phase transition that we describe above for the polar and objective function in the case of noise-free low-rank data still appears for noisy full-rank data drawn from either the independent or disjoint model as illustrated in Figure 3.7.

In order to get a better understanding of this sensitivity of the matrix factorization formulation to the generating factorization size, denoted by r_{data} , we consider the case $\gamma = 1$ (i.e. explicit column regularization on U and sparsity on V) for two reasons: first, as mentioned in the previous section, we can solve the associated polar problem in closed form; this allows us to verify that the described phenomenon is not due to approximation errors of the polar problem discussed earlier, and second, we know a (trivial) global solution to the SDL problem, $X = XI_N$, but there might exist a more compact global solution which

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

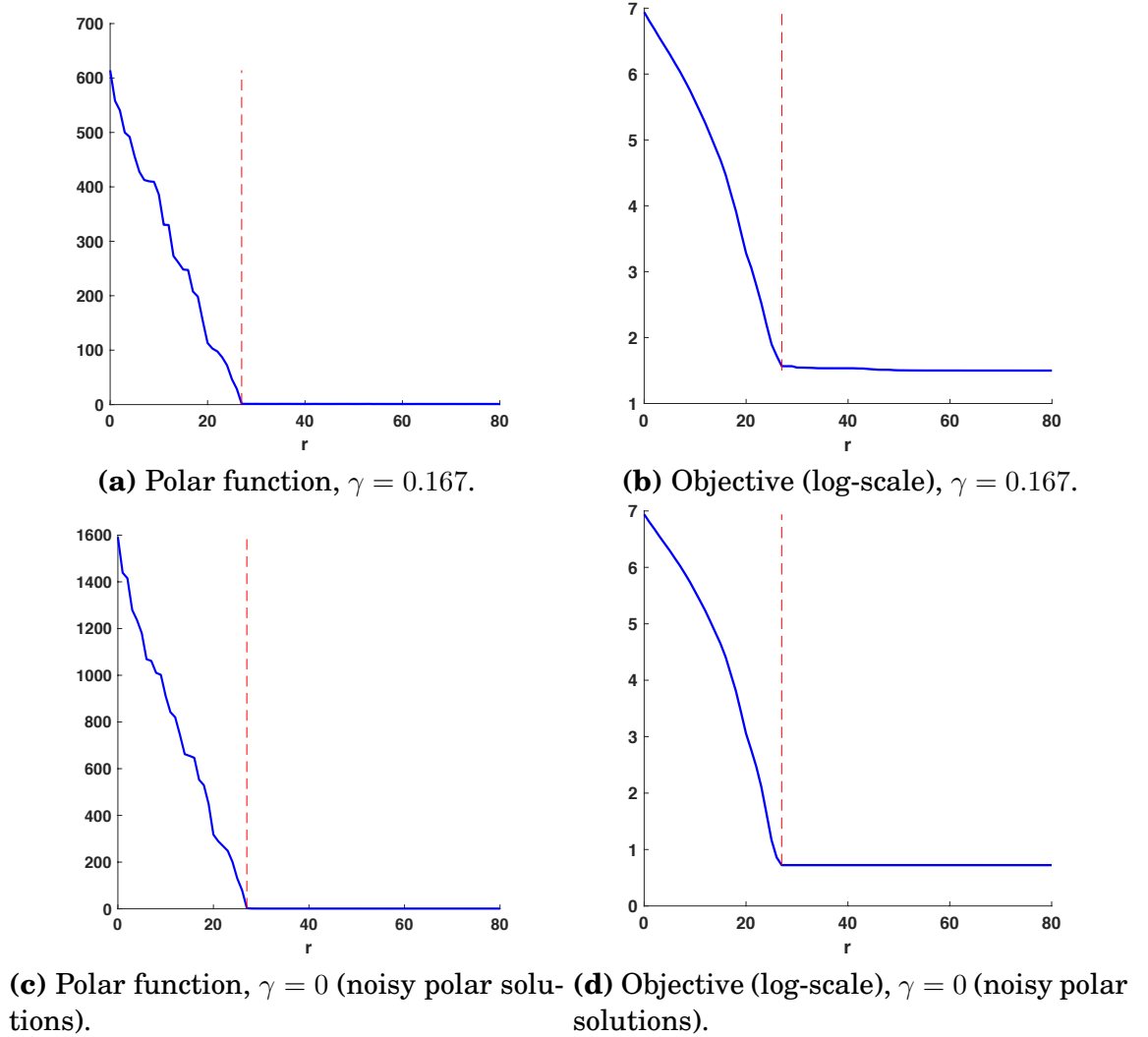


Figure 3.6: Evolution of polar and objective functions (when using approximate/noisy polar solutions). The x -axis represents the factorization size r , which coincides with the meta-iterations in this case. Problem dimensions: $X = UV^T \in \mathbb{R}^{30 \times 900}$ and $\text{rank}(X) = 27$. The red dotted line corresponds to $r_{data} = 27$.

provides more information on the actual structure of the data and which can potentially be discovered by the forward SDL MF algorithm. Intuitively such a compact solution should resemble the generating factorization, for instance in our synthetic data setting, these generating factors are U_{data} and V_{data} which we used to generate the data X . More formally, at each iteration we evaluate

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

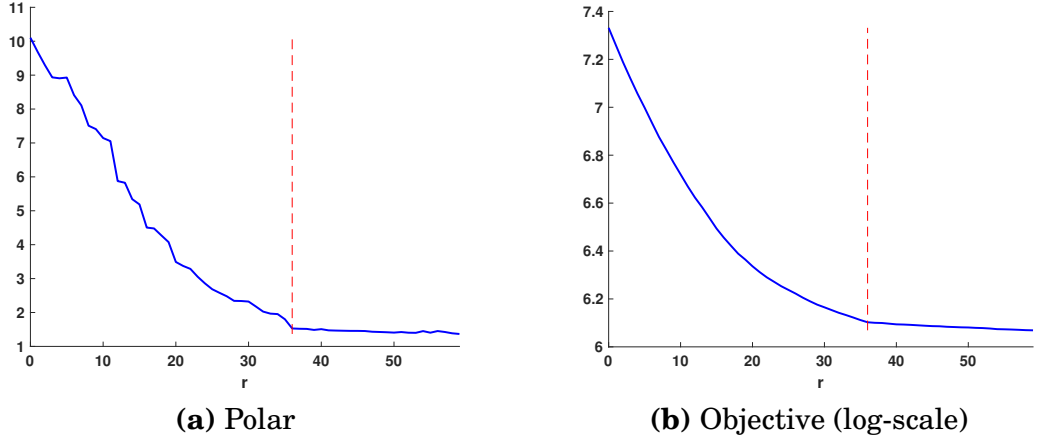


Figure 3.7: Polar and objective function throughout the meta-iterations of algorithm 1, initialized with the empty factorization. The x -axis represents the factorization size r . Problem dimensions: $X = UV^T + E \in \mathbb{R}^{30 \times 960}$, $(U, V) \in \mathbb{R}^{30 \times 36} \times \mathbb{R}^{960 \times 36}$, $E_{ij} \sim \mathcal{N}(0, 0.1^2)$. $\gamma = 0.17$. Data drawn from 12 disjoint low-dimensional subspaces of dimension 3.

$\Omega_{\theta_1} \left(-\frac{1}{\lambda} \nabla_{\Phi_k} \ell(X, \Phi_r(U^{(k)}, V^{(k)})) \right) = \Omega_{\theta_1} \left(\frac{1}{\lambda} (X - U^{(k)} V^{(k)T}) \right)$ where we recall that $\Omega_{\theta_1}(Z) = \max_i \|Z_i\|_2$. Therefore we are in fact evaluating the amplitude of the largest direction in the residual, $X - U^{(k)} V^{(k)T}$. Assume that $k < r_{data}$ and that $U^{(k)}$ is the concatenation of k basis elements of the n subspaces. Without loss of generality, assume that $U^{(k)}$ contains the bases of exactly $m < n$ subspaces S_1, \dots, S_m . Given $U^{(k)}$, trying to sparsely decompose points X_i which belong to subspaces S_{m+1}, \dots, S_n will then lead to large residuals since in non degenerate cases these points cannot be expressed as sparse linear combinations of elements from subspaces S_1, \dots, S_m .

Let us illustrate this behavior on a small example: we take U_{data} as the concatenation of 12 bases in ambient space \mathbb{R}^D (generated at random) and V_{data} as a block-diagonal matrix such that each signal comes from one of the 12 subspaces (each of dimension $d = 3$), i.e. $r_{data} = 36$. We report the results in Figure 3.8 where we observe that, up to a permutation of the columns, the factor $V^{(36)}$ (i.e. after 36 iterations) is block-diagonal and each signal is encoded by 3 atoms. This can be seen more explicitly if we compute the matrix $A = |V||V|^T$ as an

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

affinity matrix and a variant of the average *subspace-sparse recovery error*⁷: consider X_i and denote $\mathcal{J}_i = \{k \in [N] : X_i \text{ and } X_k \text{ are from the same subspace}\}$.

$$\text{ssr error} = \frac{1}{N} \sum_{i=1}^N \left(1 - \frac{\|(A_i)_{\mathcal{J}_i^c}\|_1}{\|A_i\|_1} \right)$$

This error becomes zero after r_{data} iterations of forward SDL MF, which translates as: the dictionary $U^{(k)}$ and sparse codes $V^{(k)}$, $k \geq r_{data}$, are such that points from different subspaces do not share any atoms.

In parallel we note a rapid decrease of the polar value for the first 36 iterations and a stagnation throughout the remaining iterations. A similar behavior also can be observed on the objective function which decreases much slower after the first 36 iterations. We expect this similarity since from [17] we know that the distance of the objective function at a given point (U, V) to the global optimum can be bounded by some linear function of the polar gap at this point (i.e. $\Omega_\theta \left(-\frac{1}{\lambda} \nabla_\Phi \ell(X, \Phi(U, V)) \right) - 1$). Similar results can be observed in the case of unbalanced subspace dimensions (i.e. there exist $(i, j) \in [n]^2$ such that $d_i \neq d_j$, not reported here).

Even though we know that after r_{data} iterations we have not reached the global optimum of our SDL formulation yet (the polar value is still larger than 1 in the given examples) there is empirical evidence that this forward selection strategy for factorizing some matrix X (under the union of subspaces model) allows to recover the generating factorization size in the case of independent and disjoint subspaces.

In order to detect the “knee” on the polar graph (for instance on Figures 3.7 and 3.8a), providing a trade-off between the complexity (factorization size) and the distance to the optimum (polar value), we minimize a simple objec-

⁷The subspace-sparse recovery error was introduced in the context of sparse subspace clustering where the data is used as a self-expressive dictionary and ideally every point is to be encoded exclusively by points from the subspace it belongs to. The ssr error measures the fraction of ℓ_1 -norm of the sparse representation which comes from points of other subspaces.

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

tive function linear in both parameters ($\hat{r} = \arg \min_r \alpha r + \beta \Omega_{\theta_\gamma}(Z^{(r)})$) where $Z^{(r)} = \frac{1}{\lambda}(X - U^{(r)}V^{(r)T})$ as illustrated on Figure 3.9.

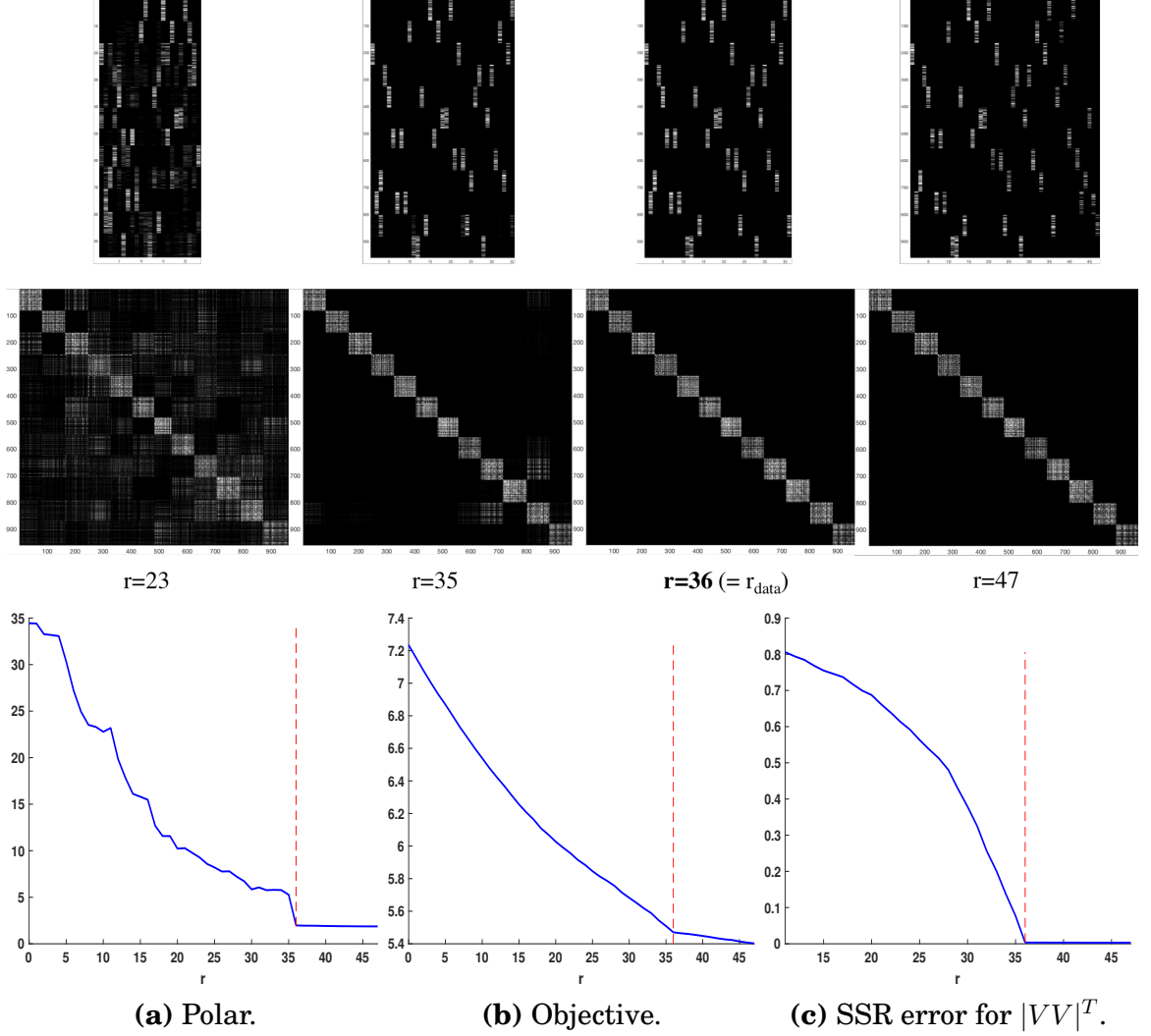


Figure 3.8: Top: Factor V (in absolute value) throughout the meta-iterations of algorithm 1, initialized with the empty factorization. Middle: Similarity matrix $|VV^T|$. Bottom: Polar and objective throughout the meta-iterations of algorithm 1. Problem dimensions: $X = UV^T \in \mathbb{R}^{30 \times 960}$, $(U, V) \in \mathbb{R}^{30 \times 36} \times \mathbb{R}^{960 \times 36}$. The data is drawn from 12 disjoint low-dimensional subspaces of dimension 3 and each entry X_{ij} is corrupted by white Gaussian noise with standard deviation $\sigma = 0.01$.

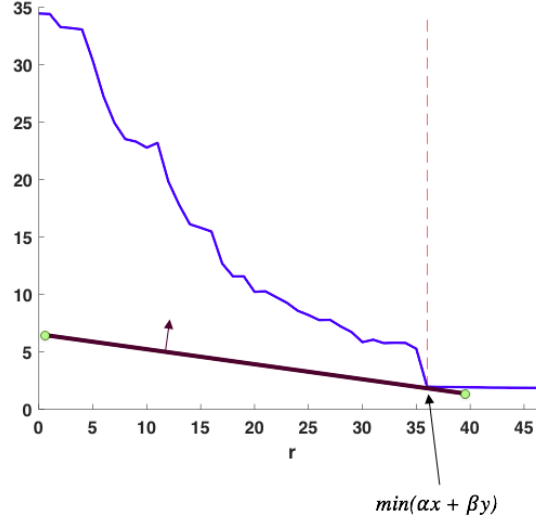


Figure 3.9: Detection of the generating factorization size on the polar curve.

3.3.3 Subspace clustering

PARAMETER TUNING

Given the previous observations and remarks we have two possibilities for choosing the sparsity parameter γ : we can take $\gamma = 1$ and follow the above described strategy to factorize the data or we tune γ on the given data since we also observed the robustness of the model selection algorithm to inexact polar solutions (which become inevitable as soon as $\gamma \in (0, 1)$). In order to tune the sparsity parameter γ which defines a trade-off between the ℓ_1 and ℓ_2 regularization on the row space, we choose a small regularization parameter $\lambda = 0.01$ to assure that the data X is a good proxy of the true solution of the convex lower bound of (2.5). We initialize a simple joint minimization (step 1 of algorithm 1) with the matrices used to construct the data matrix X ($U^{(0)} = [U_1^1 \dots U_{d_1}^1 \dots U_1^n \dots U_{d_n}^n]$) and the block-diagonal $V^{(0)} = [A_1 \dots A_N]^T$ where if $X_i \in S_j$, A_i contains the decomposition of X_i in the basis of S_j and is zero elsewhere. Due to the computational cost of the (slow) local descent we only test on a predefined set of possible values $\gamma \in \{0, 0.17, 0.33, 0.5, 0.67, 0.83, 1\}$. Assuming that the generated factors $(U^{(0)}, V^{(0)})$ are a proxy of the true so-

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

lution for a certain γ , we compare the value of the objective f at the converged result $(U^{final}, V^{final})_\gamma$ to $f(U^{(0)}, V^{(0)}) = f(U_{data}, V_{data})$ and choose $\gamma^* = \arg \min_\gamma \{ |f(U^{(0)}, V^{(0)}) - f((U^{final}, V^{final})_\gamma)| \}$.

CLUSTERING PERFORMANCE

We report the clustering performance of the forward SDL MF algorithm on data drawn from different union of subspaces models⁸. We vary the problem difficulty by generating independent and disjoint subspaces and corrupting the data by more or less noise. We compare the results to some of the previously mentioned subspace clustering methods⁹. For our method we also report the average size of the factorized solutions and their variance. For $\gamma < 1$, we choose $\gamma = 0.17$ according to the previous parameter tuning results¹⁰. As we can observe, the clustering error – defined as the fraction of points which are assigned to the wrong subspace – we achieve with our method and a tuned γ is similar to the previously presented existing methods, both for independent and disjoint subspaces. Furthermore the generating factorization size is well recovered, for both settings, as long as the noise level is small enough ($\text{SNR} \geq 4$ and $\text{SNR} \geq 10$ resp.). The failure cases with high clustering error ($\text{SNR} = 6.5$ and $\text{SNR} = 4$ in the disjoint setting) can be explained by a bad estimation of the generating factorization size r_{data} . Evaluation the clustering error at the factorization found by our algorithm after 36 iterations, i.e. using $V^{(r_{data})}$, we obtain clustering errors of 3.65% and 7.78%, respectively.

However this additionally recovered information comes at a cost: the computational efficiency of our method falls behind w.r.t. the other methods due to the many necessary local descent iterations and polar approximations.

⁸We cluster the data by Spectral Clustering on the binary k nearest neighbors affinity matrix based on $|VV|^T$.

⁹Where necessary, we tuned the parameters of these methods by grid search around the default choices suggested by the respective authors.

¹⁰The results are not very sensitive to small changes of this parameter γ , for instance for $\gamma = 0.3$ we obtain similar performances as the ones we report in Table 3.1.

CHAPTER 3. SUBSPACE CLUSTERING VIA SDL

Problem	SDL MF SC		En-SC	SSC	LR-SC
	$\gamma < 1$	$\gamma = 1$			
Independent ($n = 9, d_1 = \dots = d_n = d = 3$)					
$\sigma = 0$ ($SNR_{dB} = \infty$)	0.00 ($\bar{r} = 27,$ $\sigma_r^2 = 0$)	3.18 ($\bar{r} = 28,$ $\sigma_r^2 = 1.1$)	0.00	0.00	0.00
$\sigma = 0.01$ ($SNR_{dB} = 30$)	0.00 ($\bar{r} = 27,$ $\sigma_r^2 = 0.1$)	3.38 ($\bar{r} = 28,$ $\sigma_r^2 = 0.68$)	0.00	0.00	0.00
$\sigma = 0.1$ ($SNR_{dB} = 10$)	0.77 ($\bar{r} = 27,$ $\sigma_r^2 = 0.49$)	9.53 ($\bar{r} = 28,$ $\sigma_r^2 = 0.54$)	0.80	1.48	1.67
$\sigma = 0.15$ ($SNR_{dB} = 6.5$)	2.78 ($\bar{r} = 27,$ $\sigma_r^2 = 0.18$)	38.0 ($\bar{r} = 24,$ $\sigma_r^2 = 75$)	2.91	4.22	4.18
$\sigma = 0.2$ ($SNR_{dB} = 4$)	6.49 ($\bar{r} = 27,$ $\sigma_r^2 = 2.5$)	83.4 ($\bar{r} = 9,$ $\sigma_r^2 = 52$)	6.42	8.60	7.80
Disjoint ($n = 12, d_1 = \dots = d_n = d = 3$)					
$\sigma = 0$ ($SNR_{dB} = \infty$)	0.00 ($\bar{r} = 36,$ $\sigma_r^2 = 0$)	2.61 ($\bar{r} = 37,$ $\sigma_r^2 = 0.90$)	0.00	0.00	0.00
$\sigma = 0.01$ ($SNR_{dB} = 30$)	0.00 ($\bar{r} = 36,$ $\sigma_r^2 = 0.1$)	2.55 ($\bar{r} = 37,$ $\sigma_r^2 = 0.71$)	0.00	0.00	0.00
$\sigma = 0.1$ ($SNR_{dB} = 10$)	1.06 ($\bar{r} = 36,$ $\sigma_r^2 = 1.2$)	9.86 ($\bar{r} = 37,$ $\sigma_r^2 = 1.1$)	1.03	1.70	2.07
$\sigma = 0.15$ ($SNR_{dB} = 6.5$)	9.63 ($\bar{r} = 27,$ $\sigma_r^2 = 6.9$)	40.6 ($\bar{r} = 27,$ $\sigma_r^2 = 209$)	3.23	4.80	4.92
$\sigma = 0.2$ ($SNR_{dB} = 4$)	17.9 ($\bar{r} = 28,$ $\sigma_r^2 = 3.6$)	96.6 ($\bar{r} = 7,$ $\sigma_r^2 = 12$)	7.44	9.62	8.90

Table 3.1: Clustering error (in %) and detected factorization size (average \bar{r} and variance σ_r^2 are taken over 10 trials).

4 Conclusion

In this work, we presented the problem of sparse dictionary learning framed in a specific structured matrix factorization formulation which allowed us to derive conditions of global optimality for the dictionary and the sparse coding variables. We studied these conditions and means to verify them by establishing a better understanding of the associated polar problem, empirically and theoretically. From empirical observations we conjecture that the number of local optima of this polar problem is linear in the dimension of the data, opening the perspective of localizing the global optimum despite the non-convexity of the problem, by optimizing over multiple initializations. A complete theoretical analysis of the optimization landscape of this polar problem is left for future work.

In order to apply the theoretical results on global optimality we proposed an iterative algorithm to find globally optimal solutions for the sparse dictionary learning problem. We motivated the use of our formulation for subspace clustering with the additional aim of finding an optimal subspace-preserving factorization of the data. This method has similar clustering performance than current state-of-the-art methods but its computational complexity discourages its direct application to subspace clustering. However it allows to recover a good estimation of the generating factorization size of the data and therefore provides a tool to identify a (small) factorization size adapted to the data model and allowing for a compact and potentially more robust representation of the data. A next step will consist in providing a theoretical analysis for this sensitivity of the polar function to the generating factorization size under adequate assumptions on the data.

A Appendix

In the following we provide the proof for Proposition 3:

Proof. We recall the optimization problem we are trying to solve:

$$\min_{\mathbf{v} \in \mathbb{R}^N} f(\mathbf{v}) := \frac{\gamma \|\mathbf{v}\|_1 + (1 - \gamma)}{\|Z\mathbf{v}\|_2} \quad \text{subject to } \|\mathbf{v}\|_2 = 1 \quad (\text{A.1})$$

First we compute the Euclidean and the Riemannian subdifferential at \mathbf{v} of the objective function f , resp. $\partial f(\mathbf{v})$ and $\text{grad } f(\mathbf{v})$:

$$\begin{aligned} \partial f(\mathbf{v}) &= \frac{\gamma \|Z\mathbf{v}\|_2 \partial \|\mathbf{v}\|_1 - \frac{(\gamma \|\mathbf{v}\|_1 + (1 - \gamma)) Z^T Z \mathbf{v}}{\|Z\mathbf{v}\|_2}}{\|Z\mathbf{v}\|_2^2} \\ &= \frac{\gamma \|Z\mathbf{v}\|_2^2 \partial \|\mathbf{v}\|_1 - (\gamma \|\mathbf{v}\|_1 + (1 - \gamma)) Z^T Z \mathbf{v}}{\|Z\mathbf{v}\|_2^3} \\ &= \frac{\gamma \eta(\mathbf{v}) \partial \|\mathbf{v}\|_1 - \kappa(\mathbf{v}) H \mathbf{v}}{\eta(\mathbf{v})^{3/2}} \end{aligned} \quad (\text{A.2})$$

and

$$\text{grad } f(\mathbf{v}) = (\mathbf{I} - \mathbf{v} \mathbf{v}^T) \partial f(\mathbf{v}) = \frac{\eta(\mathbf{v}) (\gamma \partial \|\mathbf{v}\|_1 + (1 - \gamma) \mathbf{v}) - \kappa(\mathbf{v}) H \mathbf{v}}{\eta(\mathbf{v})^{3/2}} \quad (\text{A.3})$$

A first-order necessary condition for A.1 is

$$0 \in \text{grad } f(\mathbf{v}) \quad \Leftrightarrow 0 \in \eta(\mathbf{v}) (\gamma \partial \|\mathbf{v}\|_1 + (1 - \gamma) \mathbf{v}) - \kappa(\mathbf{v}) H \mathbf{v} \quad (\text{A.4})$$

Partitioning this condition according to the support of \mathbf{v} , σ and noticing that $H_{11} \mathbf{v}_1 = H_{11} S_1 S_1^T \mathbf{v}_1 = H_{11} S_1 \tilde{\mathbf{v}}_1$, we obtain the restricted first-order optimality conditions FO_σ and FO_{σ^c} .

In order to derive the second-order necessary condition we take problem A.1

APPENDIX A. APPENDIX

restricted to the sign of v :

$$\min_{\mathbf{w} \in \mathbb{R}^N} f(\mathbf{w}) \quad \text{subject to } \|\mathbf{w}\|_2 = 1, \text{ supp}(\mathbf{w}) \subseteq \sigma, S_1 \mathbf{w}_\sigma \geq \mathbf{0} \quad (\text{A.5})$$

We can further obtain a smooth optimization problem over the manifold $\mathbb{S}^{|\sigma|-1} \cap \mathbb{R}_+^{|\sigma|}$ by using the change of variables $\mathbf{w} : \mathbf{x} \mapsto \mathbf{w}(\mathbf{x}) = [S_1 \quad \mathbf{0}]^T \mathbf{x}$.

$$\min_{\mathbf{x} \in \mathbb{R}^{|\sigma|}} \bar{f}(\mathbf{x}) \quad \text{subject to } \|\mathbf{x}\|_2 = 1, \mathbf{x} \geq \mathbf{0} \quad (\text{A.6})$$

We will derive the second-order necessary condition for A.6 since if \tilde{v}_1 satisfies this condition, then v is a local minimum of A.1. The second-order necessary condition for a smooth objective \bar{f} over the sphere is given in [1, 2] as

$$\langle \xi, \text{Hess } \bar{f}(\mathbf{x})[\xi] \rangle = \xi^T \nabla^2 \bar{f}(\mathbf{x}) \xi - (\mathbf{x}^T \nabla \bar{f}(\mathbf{x})) \|\xi\|_2^2 \geq 0 \quad \text{for all } \xi \in \mathcal{T}_{\mathbf{x}} \mathbb{S}^{|\sigma|-1} \quad (\text{A.7})$$

where we denote by $\text{Hess } \bar{f}(\mathbf{x})$ and $\nabla^2 \bar{f}(\mathbf{x})$ the Riemannian and Euclidian Hessian respectively.

Similar to the computation of $\nabla f(v)$, we obtain the gradient of $\bar{f}(\tilde{v}_1)$ as

$$\nabla \bar{f}(\tilde{v}_1) = \frac{\gamma \eta(\tilde{v}_1) \mathbf{1} - \kappa(\tilde{v}_1) \tilde{H}_{11} \tilde{v}_1}{\eta(\tilde{v}_1)^{3/2}} \quad (\text{A.8})$$

APPENDIX A. APPENDIX

And the Hessian $\nabla^2 \bar{f}(\tilde{\mathbf{v}}_1)$ as

$$\begin{aligned}
\nabla^2 \bar{f}(\tilde{\mathbf{v}}_1) &= (\gamma\eta(\tilde{\mathbf{v}}_1)\mathbf{1} - \kappa(\tilde{\mathbf{v}}_1)\tilde{H}_{11}\tilde{\mathbf{v}}_1) (\nabla(\eta(\tilde{\mathbf{v}}_1))^{-1})^T + (\eta(\tilde{\mathbf{v}}_1))^{-1} \mathbf{J}(\gamma\eta(\tilde{\mathbf{v}}_1)\mathbf{1} - \kappa(\tilde{\mathbf{v}}_1)\tilde{H}_{11}\tilde{\mathbf{v}}_1) \\
&= \frac{2\gamma\mathbf{1}(\tilde{H}_{11}\tilde{\mathbf{v}}_1)^T - \gamma\tilde{H}_{11}\tilde{\mathbf{v}}_1\mathbf{1}^T - (\gamma\mathbf{1}^T\tilde{\mathbf{v}}_1 + (1-\gamma))\tilde{H}_{11}}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} \\
&\quad - 3 \frac{\left(\gamma\eta(\tilde{\mathbf{v}}_1)\mathbf{1} - (\gamma\mathbf{1}^T\tilde{\mathbf{v}}_1 + (1-\gamma))\tilde{H}_{11}\tilde{\mathbf{v}}_1 \right) (\tilde{H}_{11}\tilde{\mathbf{v}}_1)^T}{\eta(\tilde{\mathbf{v}}_1)^{5/2}} \\
&= \frac{2\gamma\mathbf{1}(\tilde{H}_{11}\tilde{\mathbf{v}}_1)^T - \gamma\tilde{H}_{11}\tilde{\mathbf{v}}_1\mathbf{1}^T - \kappa(\tilde{\mathbf{v}}_1)\tilde{H}_{11}}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} - 3 \frac{\left(\gamma\eta(\tilde{\mathbf{v}}_1)\mathbf{1} - \kappa(\tilde{\mathbf{v}}_1)\tilde{H}_{11}\tilde{\mathbf{v}}_1 \right) (\tilde{H}_{11}\tilde{\mathbf{v}}_1)^T}{\eta(\tilde{\mathbf{v}}_1)^{5/2}} \\
&= \frac{3\kappa(\tilde{\mathbf{v}}_1)}{\eta(\tilde{\mathbf{v}}_1)^{5/2}} (\tilde{H}_{11}\tilde{\mathbf{v}}_1)(\tilde{H}_{11}\tilde{\mathbf{v}}_1)^T - \frac{\gamma}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} (\mathbf{1}(\tilde{H}_{11}\tilde{\mathbf{v}}_1)^T + (\tilde{H}_{11}\tilde{\mathbf{v}}_1)\mathbf{1}^T) - \frac{\kappa(\tilde{\mathbf{v}}_1)}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} \tilde{H}_{11}
\end{aligned} \tag{A.9}$$

where $\mathbf{J}(g)$ is the Jacobian matrix of some vector-valued function g .

Assuming \mathbf{v} satisfies FO_σ , we can simplify this expression to obtain

$$\begin{aligned}
\nabla^2 \bar{f}(\tilde{\mathbf{v}}_1) &= \frac{3\kappa(\tilde{\mathbf{v}}_1)}{\eta(\tilde{\mathbf{v}}_1)^{5/2}} \left(\frac{\eta(\tilde{\mathbf{v}}_1)}{\kappa(\tilde{\mathbf{v}}_1)} \right)^2 (\gamma\mathbf{1} + (1-\gamma)\tilde{\mathbf{v}}_1)(\gamma\mathbf{1} + (1-\gamma)\tilde{\mathbf{v}}_1)^T \\
&\quad - \frac{\gamma}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} \frac{\eta(\tilde{\mathbf{v}}_1)}{\kappa(\tilde{\mathbf{v}}_1)} (\mathbf{1}(\gamma\mathbf{1} + (1-\gamma)\tilde{\mathbf{v}}_1)^T + (\gamma\mathbf{1} + (1-\gamma)\tilde{\mathbf{v}}_1)\mathbf{1}^T) \\
&\quad - \frac{\kappa(\tilde{\mathbf{v}}_1)}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} \tilde{H}_{11} \\
&= \frac{\gamma^2}{\eta(\tilde{\mathbf{v}}_1)^{1/2}\kappa(\tilde{\mathbf{v}}_1)} \mathbf{1}\mathbf{1}^T + 3 \frac{(1-\gamma)^2}{\eta(\tilde{\mathbf{v}}_1)^{1/2}\kappa(\tilde{\mathbf{v}}_1)} \tilde{\mathbf{v}}_1\tilde{\mathbf{v}}_1^T \\
&\quad + 2 \frac{\gamma(1-\gamma)}{\eta(\tilde{\mathbf{v}}_1)^{1/2}\kappa(\tilde{\mathbf{v}}_1)} (\mathbf{1}\tilde{\mathbf{v}}_1^T + \tilde{\mathbf{v}}_1\mathbf{1}^T) - \frac{\kappa(\tilde{\mathbf{v}}_1)}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} \tilde{H}_{11}
\end{aligned} \tag{A.10}$$

We integrate this result into the second-order optimality condition A.7 to

APPENDIX A. APPENDIX

obtain

$$\begin{aligned} & \xi^T \left(\frac{\gamma^2}{\eta(\tilde{\mathbf{v}}_1)^{1/2} \kappa(\tilde{\mathbf{v}}_1)} \mathbf{1}\mathbf{1}^T + 3 \frac{(1-\gamma)^2}{\eta(\tilde{\mathbf{v}}_1)^{1/2} \kappa(\tilde{\mathbf{v}}_1)} \tilde{\mathbf{v}}_1 \tilde{\mathbf{v}}_1^T \right) \xi \\ & + \xi^T \left(2 \frac{\gamma(1-\gamma)}{\eta(\tilde{\mathbf{v}}_1)^{1/2} \kappa(\tilde{\mathbf{v}}_1)} (\mathbf{1} \tilde{\mathbf{v}}_1^t + \tilde{\mathbf{v}}_1 \mathbf{1}^T) - \frac{\kappa(\tilde{\mathbf{v}}_1)}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} \tilde{H}_{11} \right) \xi \end{aligned} \quad (\text{A.11})$$

$$\begin{aligned} & - (\tilde{\mathbf{v}}_1^T \nabla \bar{f}(\tilde{\mathbf{v}}_1)) \|\xi\|_2^2 \geq 0 \quad \text{for all } \xi \in \mathcal{T}_x \mathbb{S}^{|\sigma|-1} \\ \Leftrightarrow & \xi^T \left(\frac{\gamma^2}{\eta(\tilde{\mathbf{v}}_1)^{1/2} \kappa(\tilde{\mathbf{v}}_1)} \mathbf{1}\mathbf{1}^T - \frac{\kappa(\tilde{\mathbf{v}}_1)}{\eta(\tilde{\mathbf{v}}_1)^{3/2}} \tilde{H}_{11} \right) \xi \\ & + \frac{1-\gamma}{\eta(\tilde{\mathbf{v}}_1)^{1/2}} \|\xi\|_2^2 \geq 0 \quad \text{for all } \xi \in \mathcal{T}_x \mathbb{S}^{|\sigma|-1} \end{aligned} \quad (\text{A.12})$$

$$\Leftrightarrow \frac{\gamma^2}{\kappa(\tilde{\mathbf{v}}_1)^2} (\mathbf{1}^T \xi)^2 + \frac{1-\gamma}{\kappa(\tilde{\mathbf{v}}_1)^2} \geq \frac{\xi^T \tilde{H}_{11} \xi}{\eta(\tilde{\mathbf{v}}_1)} \quad \text{for all } \xi \in \mathcal{T}_x \mathbb{S}^{|\sigma|-1} \quad (\text{A.13})$$

□

Bibliography

- [1] ABSIL, P.-A., MAHONY, R., AND SEPULCHRE, R. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- [2] ABSIL, P.-A., MAHONY, R., AND TRUMPF, J. An extrinsic look at the riemannian hessian. In *Geometric science of information*. Springer, 2013, pp. 361–368.
- [3] ADLER, A., ELAD, M., AND HEL-OR, Y. Linear-time subspace clustering via bipartite graph modeling. *IEEE transactions on neural networks and learning systems* 26, 10 (2015), 2234–2246.
- [4] BACH, F. Convex relaxations of structured matrix factorizations. *arXiv preprint arXiv:1309.3117* (2013).
- [5] BACH, F., MAIRAL, J., AND PONCE, J. Convex sparse matrix factorizations. *arXiv preprint arXiv:0812.1869* (2008).
- [6] CAI, J.-F., CANDÈS, E. J., AND SHEN, Z. A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization* 20, 4 (2010), 1956–1982.
- [7] CAVAZZA, J., MORERIO, P., HAEFFELE, B., LANE, C., MURINO, V., AND VIDAL, R. Dropout as a low-rank regularizer for matrix factorization. *arXiv preprint arXiv:1710.05092* (2017).
- [8] CHEN, Y., AND CANDÈS, E. J. The projected power method: An efficient algorithm for joint alignment from pairwise differences. *Communications on Pure and Applied Mathematics* 71, 8 (2018), 1648–1714.
- [9] EFRON, B., HASTIE, T., JOHNSTONE, I., TIBSHIRANI, R., ET AL. Least angle regression. *The Annals of statistics* 32, 2 (2004), 407–499.

BIBLIOGRAPHY

- [10] ELAD, M., AND AHARON, M. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image processing* 15, 12 (2006), 3736–3745.
- [11] ELHAMIFAR, E., AND VIDAL, R. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE transactions on pattern analysis and machine intelligence* 35, 11 (2013), 2765–2781.
- [12] FAVARO, P., VIDAL, R., AND RAVICHANDRAN, A. A closed form solution to robust subspace estimation and clustering. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on* (2011), IEEE, pp. 1801–1807.
- [13] GE, R., HUANG, F., JIN, C., AND YUAN, Y. Escaping from saddle points – online stochastic gradient for tensor decomposition. In *Conference on Learning Theory* (2015), pp. 797–842.
- [14] GE, R., AND MA, T. On the optimization landscape of tensor decompositions. In *Advances in Neural Information Processing Systems* (2017), pp. 3653–3663.
- [15] GRANT, M., BOYD, S., AND YE, Y. Cvx: Matlab software for disciplined convex programming, 2008.
- [16] GUNASEKAR, S., LEE, J., SOUDRY, D., AND SREBRO, N. Characterizing implicit bias in terms of optimization geometry. *arXiv preprint arXiv:1802.08246* (2018).
- [17] HAEFFELE, B. D., AND VIDAL, R. Global optimality in tensor factorization, deep learning, and beyond. *arXiv preprint arXiv:1506.07540* (2015).
- [18] HAEFFELE, B. D., AND VIDAL, R. Structured low-rank matrix factorization: Global optimality, algorithms, and applications. *arXiv preprint arXiv:1708.07850* (2017).

BIBLIOGRAPHY

- [19] HEIN, M., AND BÜHLER, T. An inverse power method for nonlinear eigenproblems with applications in 1-spectral clustering and sparse pca. In *Advances in Neural Information Processing Systems* (2010), pp. 847–855.
- [20] HENDRICKX, J. M., AND OLSHEVSKY, A. Matrix p-norms are np-hard to approximate if $p \neq 1, 2, \infty$. *SIAM Journal on Matrix Analysis and Applications* 31, 5 (2010), 2802–2812.
- [21] JOURNÉE, M., NESTEROV, Y., RICHTÁRIK, P., AND SEPULCHRE, R. Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research* 11, Feb (2010), 517–553.
- [22] LEE, D. D., AND SEUNG, H. S. Learning the parts of objects by non-negative matrix factorization. *Nature* 401, 6755 (1999), 788.
- [23] LIU, G., LIN, Z., AND YU, Y. Robust subspace segmentation by low-rank representation. In *Proceedings of the 27th international conference on machine learning (ICML-10)* (2010), pp. 663–670.
- [24] MAIRAL, J., BACH, F., PONCE, J., ET AL. Sparse modeling for image and vision processing. *Foundations and Trends® in Computer Graphics and Vision* 8, 2-3 (2014), 85–283.
- [25] MAIRAL, J., BACH, F., PONCE, J., AND SAPIRO, G. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research* 11, Jan (2010), 19–60.
- [26] MAIRAL, J., PONCE, J., SAPIRO, G., ZISSERMAN, A., AND BACH, F. R. Supervised dictionary learning. In *Advances in neural information processing systems* (2009), pp. 1033–1040.
- [27] MALLAT, S. *A wavelet tour of signal processing: the sparse way*. Academic press, 2008.

BIBLIOGRAPHY

- [28] OLSHAUSEN, B. A., AND FIELD, D. J. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research* 37, 23 (1997), 3311–3325.
- [29] PARIKH, N., BOYD, S., ET AL. Proximal algorithms. *Foundations and Trends® in Optimization* 1, 3 (2014), 127–239.
- [30] RECHT, B., FAZEL, M., AND PARRILO, P. A. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review* 52, 3 (2010), 471–501.
- [31] RODRIGUEZ, F., AND SAPIRO, G. Sparse representations for image classification: Learning discriminative and reconstructive non-parametric dictionaries. Tech. rep., University of Minnesota. Institute for Mathematics and Its Applications, 2008.
- [32] SCHWAB, E., HAEFFELE, B., CHARON, N., AND VIDAL, R. Separable dictionary learning with global optimality and applications to diffusion mri. *arXiv preprint arXiv:1807.05595* (2018).
- [33] SEFATI, S., COWAN, N. J., AND VIDAL, R. Learning shared, discriminative dictionaries for surgical gesture segmentation and classification. In *MICCAI Workshop: M2CAI* (2015), vol. 4.
- [34] SUN, J., QU, Q., AND WRIGHT, J. Complete dictionary recovery over the sphere. *arXiv preprint arXiv:1504.06785* (2015).
- [35] TIBSHIRANI, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* (1996), 267–288.
- [36] VIDAL, R., MA, Y., AND SASTRY, S. S. *Generalized principal component analysis*, vol. 5. Springer, 2016.
- [37] WRIGHT, J., YANG, A. Y., GANESH, A., SASTRY, S. S., AND MA, Y. Robust face recognition via sparse representation. *IEEE transactions on pattern analysis and machine intelligence* 31, 2 (2009), 210–227.

BIBLIOGRAPHY

- [38] ZHU, Z., LI, Q., TANG, G., AND WAKIN, M. B. Global optimality in low-rank matrix optimization. *IEEE Transactions on Signal Processing* 66, 13 (2018), 3614–3628.
- [39] ZOU, H., HASTIE, T., AND TIBSHIRANI, R. Sparse principal component analysis. *Journal of computational and graphical statistics* 15, 2 (2006), 265–286.