

**Analysis and Extension of Sparse Representations in Signal Classification**

By

CHELSEA A. WEAVER

B.S. (University of Washington) 2009

DISSERTATION

Submitted in partial satisfaction of the requirements for the degree of

DOCTOR OF PHILOSOPHY

in

MATHEMATICS

in the

OFFICE OF GRADUATE STUDIES

of the

UNIVERSITY OF CALIFORNIA

DAVIS

Approved:

---

Naoki Saito (Chair)

---

Thomas Strohmer

---

Jesus De Loera

Committee in Charge

2016

© Chelsea A. Weaver, 2016. All rights reserved.

To Chris, for firmly telling me to stay the course. Consider the river crossed.

# Contents

Abstract	vi
Acknowledgments	viii
Chapter 1. Introduction	1
1.1. Notation	5
1.2. Lists of Figures and Tables	6
<b>Part 1. Improving Representation-Based Classification</b>	<b>10</b>
Chapter 2. Classification	11
2.1. The Task	11
2.2. Pre-Processing and Feature Extraction	15
2.3. The Class Manifold Assumption	17
Chapter 3. Representation-Based Classification	19
3.1. Per-Class Decomposition	19
3.2. Sparse Representation-Based Classification	22
3.3. Similar Methods	25
3.4. Summary	32
Chapter 4. Local Principal Component Analysis SRC	33
4.1. Motivation	33
4.2. Algorithm Statement	34
4.3. Local Principal Component Analysis	35
4.4. Remarks on the Choice of Parameters	36
4.5. Computational Complexity and Storage Requirements	40

Chapter 5. Experiments with LPCA-SRC	42
5.1. Algorithms Compared	42
5.2. Setting of Parameters	43
5.3. Synthetic Database	44
5.4. Face Databases	50
Chapter 6. Bounding the Tangent Error	61
6.1. Description and Assumptions	61
6.2. Kaslovsky and Meyer's Tangent Bound	62
6.3. Main Theorem	63
6.4. Proof of Main Theorem	64
6.5. Remarks Regarding Implementation	75
Chapter 7. Other Local PCA Modifications	78
7.1. Modifying CRC-RLS	78
7.2. Modifying Structured Sparsity/Block-Sparse Methods	84
7.3. Discussion	91
Chapter 8. Conclusion of Part 1	93
<b>Part 2. Examining Sparsity in Classification</b>	97
Chapter 9. Equivalence Guarantees	98
9.1. Motivation from Compressed Sensing	98
9.2. Review of Equivalence Guarantees	100
Chapter 10. Mutual Coherence Equivalence in the Context of Classification	107
10.1. The Dilemma	108
10.2. Main Goal	109
10.3. Preliminary Results	109
10.4. Main Result	110
Chapter 11. Equivalence on Highly-Coherent Data	115

11.1.	Inspiration	115
11.2.	Project Description	117
11.3.	Experiments	118
11.4.	Summary	130
Chapter 12. Proving Equivalence in SRC via Nonlinear Embedding		132
12.1.	This Approach	132
12.2.	Designing the Transform	134
12.3.	An Unfortunate yet Necessary Modification	138
12.4.	Experiments	141
12.5.	Key Findings	146
Chapter 13. Conclusion of Part 2		151
Chapter 14. Final Remarks		153
Appendix A. Proof of Tangent Bound Modification		156
A.1.	Tangent Bound Details	156
A.2.	Singer and Wu's Local PCA	159
A.3.	Modifying Kaslovsky and Meyer's Setup	160
A.4.	Scaling Factor	162
A.5.	Bounding Terms	165
A.6.	End Result	184
A.7.	Adjusting the Scaling Factor	188
Bibliography		190

Chelsea A. Weaver  
August 2016  
Mathematics

## Analysis and Extension of Sparse Representations in Signal Classification

### Abstract

Classification, the task of assigning class labels to incoming data samples based on a training set, has as many applications as it has approaches. Our focus here is on classification algorithms that seek to represent the given test sample as an (often regularized) linear combination of the training samples. This dissertation has two parts:

First, we assess the representation-based approach to classification and the assumptions required for such methods to be effective. Wright et al.'s popular *sparse representation-based classification* (SRC) will be at the forefront of our analysis. By minimizing the  $\ell^1$ -norm of the coefficient vector, SRC seeks a *sparse* decomposition of the test sample over the set of training samples, with classification to the most-contributing class. Though this algorithm is largely successful, it assumes that class manifolds are linear subspaces spanned by their training samples, which is often impractical. We proceed to modify SRC so that this restrictive assumption may be relaxed and the applicability of the algorithm extended. In particular, we use local principal component analysis to approximate the tangent hyperplane of the class manifold at each training sample, and we then use the basis vectors of these tangent hyperplanes to supplement a carefully-selected subset of training samples over which the test sample is sparsely decomposed. We show that the resulting, novel algorithm leads to better classification rates than SRC in cases of sparsely-sampled and/or nonlinear class manifolds, low noise, and stringent dimensionality reduction. Additionally, we state and prove a theorem that bounds the distance between the tangent vectors computed in our algorithm and the (ground truth) class manifolds. We also validate the use of sparsity (via  $\ell^1$ -minimization) in the representation framework by comparing our modification of SRC to similar modifications of (i)  $\ell^2$ -regularized (instead of  $\ell^1$ -regularized) representation-based classification and (ii) a structured sparsity/block-sparse (minimization of the number of *classes* whose training samples have nonzero coefficients) approach to representation-based classification.

Second, we investigate the relationship between classification and sparsity, specifically, the two-fold question of whether sparsity in the representation is necessary for good classification performance, and whether or not the sparsest solution can be efficiently obtained in practice. Though under certain conditions, minimization of the  $\ell^1$ -norm provably recovers the sparsest solution, we determine that the tractable approach to verifying this sparse recovery is fundamentally in conflict with the high coherence between training samples in the same class. Despite the lack of implementable recovery guarantee, we show experimentally that the sparsest solution is often recovered by  $\ell^1$ -minimization in the case that the classes are well-separated. Further, through the use of a nonlinear transform designed so that sparse recovery conditions may be satisfied, we examine the relationship between the accuracy of SRC and the sparsity of its representation. We demonstrate that *approximate* (and not strict) equivalence between the  $\ell^1$ -minimized solution and the sparsest solution is key to the success of SRC.

## Acknowledgments

My research was conducted with government support under contract FA9550-11-C-0028 and awarded by DoD, Air Force Office of Scientific Research, National Defense Science and Engineering Graduate (NDSEG) Fellowship, 32 CFR 168a, as well as by National Science Foundation VIGRE DMS-0636297 and NSF DMS-1418779. My adviser, Naoki Saito, was supported by ONR grants N00014-12-1-0177 and N00014-16-1-2255, as well as NSF DMS-1418779, during my time at UC Davis.

I would like to acknowledge the help of François Meyer and Daniel Kaslovsky in understanding parts of their research. They responded quickly and very kindly to my emails, and their help greatly surpassed my expectations and aided the tangent bound modification in Appendix A of this dissertation.

My sincere thanks to Naoki Saito for his meticulous feedback and supportive mentoring style, as well as for taking me on as his student before I had any real idea of what having an adviser meant. I'd also like to thank my committee members, Thomas Strohmer and Jesus De Loera—your encouragement and belief in me means more than I can say. Additionally, I could not have done this without my big-office peers Ryan Halabi and Jeff Irion, for whom I had more questions than answers, and Calina Copos, who literally made the filing of this dissertation possible.

Lastly, I'd like to profusely thank my family and friends for all their patience, encouragement, and support. They had unwavering faith in me every step of the way and gave me praise and forgiveness when I needed it most.

## CHAPTER 1

# Introduction

*Classification* is the task of assigning labels to test samples (samples whose classes are unknown) given the class information of a training set. Many classification algorithms have been proposed for various important applications, including detection of liver cancer from CT scans in medical image analysis [76], loan default prediction via credit scoring [47], earthquake detection [6], and human action recognition [66]. The classification of handwritten digits has been applied to the computer-automated reading of postal zip codes [58], and algorithms for face recognition, critical to applications from national security to social media, have been proposed [17, 71, 104]. As we will see, there are many criteria involved in assessing the goodness of a classification method, and many proposed classifiers offer an improvement in only some specific area(s). Along with the deluge of big data and recent advances in artificial intelligence, the difficulty of the task drives the seemingly unending need for better (faster, more accurate, more informative) classification algorithms.

In this dissertation, we focus on classification algorithms that seek to *express* the given test sample using the set of training samples. It follows that training samples that occur prominently in the expression (e.g., representation or approximation) are in some way deeply similar to the test sample, and so we can correctly infer the test sample's class from the class(es) of these significant training samples. For instance, when the test sample is expressed as a linear combination of the training samples, the (absolute values of the) coefficients can indicate the weight or degree of similarity between the test sample and each training sample. Thus we can pick out significant training samples by looking at their coefficients. We will refer to this as *representation-based classification*. This dissertation has two parts:

In Part 1, we assess the representation-based approach to classification and the assumptions required for such methods to be effective. In particular, we focus on algorithms that express the test sample as a regularized linear combination of the entire set of training samples. The popular algorithm *sparse representation-based classification* (SRC) [104] will be at the forefront of our

---

analysis. Proposed in 2009 by Wright et al., SRC was motivated by the recent boom in the use of sparse representation in signal processing. The catalyst of these advancements was the discovery that the sparsest representation of a signal using an over-complete set of vectors (often called a *dictionary*) could be found in many situations by minimizing the  $\ell^1$ -norm of the representation coefficient vector [29]. Since the  $\ell^1$ -minimization problem is convex, this gave rise to a tractable approach to obtaining the sparsest solution in many signal processing applications.

SRC applies this relationship between the minimum  $\ell^1$ -norm and the sparsest solution to classification. The algorithm seeks the sparsest decomposition of a test sample over the dictionary of training samples via  $\ell^1$ -minimization, with classification to the class whose corresponding portion of the representation approximates the test sample with least error. The method assumes that class manifolds are linear subspaces, so that the test sample can be represented using training samples in its ground truth class. Wright et al. [104] argue that this is precisely the sparsest decomposition of the test sample over the training set. They make the case that sparsity is critical to high-dimensional image classification and that, if properly harnessed, it can lead to impressive classification performance, even on highly-corrupted or occluded images. Further, good results can be achieved regardless of the choice of image features that are used for classification, provided that the number of retained features is large enough [104]. Though SRC was originally applied to face recognition, similar methods have been employed in clustering [39], dimensionality reduction [73], and texture and handwritten digit classification [108].

As we will see, the effectiveness of SRC and similar representation-based algorithms hinges on the lower-dimensional structure of each class manifold and the size and representativeness of the training set. When the data is poorly-sampled from nonlinear manifolds, it may be impossible to express test samples as approximate linear combinations of their same-class training samples.

We proceed to modify SRC to specifically address these issues. We call the resulting algorithm *local principal component analysis sparse representation-based classification* (LPCA-SRC). Our algorithm is designed to increase the training set via basis vectors of the hyperplanes approximately tangent to the (unknown) class manifolds. This provides the two-fold benefit of counter-balancing the potential sparse sampling of class manifolds (especially in the case that they are nonlinear) and helping to regain some information lost due to the dimensionality reduction procedure applied to

---

the data set as pre-processing. The dictionary of training samples in SRC is replaced with a set of vectors consisting of training samples that are close to the given test sample and their corresponding tangent hyperplane basis vectors, increasing the accuracy and locality of the approximation of the test sample in terms of its ground truth class. We demonstrate that this modified dictionary leads to classification accuracy exceeding that of traditional SRC and related methods on a synthetic database and three popular face databases.

Towards quantifying the efficacy of LPCA-SRC, we next use a result by Kaslovsky and Meyer [57] to bound the maximum distance between the tangent vectors computed in LPCA-SRC and the ground truth class manifolds. This analysis produces a potential tool for the user to gauge whether our modification of SRC can improve classification accuracy (over the original SRC) on a given data set, under the framework that the produced tangent vectors are newly-generated, potentially noisy training samples.

We conclude Part 1 by further validating the importance of the sparse-representation framework in LPCA-SRC (i.e., the “SRC” in “LPCA-SRC”). In particular, we similarly modify two other representation-based classification algorithms (a method based on  $\ell^2$ -regularization and one based on structured/block-sparsity) using tangent vectors, and we show that our modification improves each algorithm but that  $\ell^1$ -minimization is needed for good classification accuracy in several cases.

In Part 2, we investigate the role of sparsity in SRC, specifically, the two-fold question of (i) whether sparsity (achieved by minimizing the number of training samples with nonzero coefficients in the representation of the test sample) can improve classification performance, and (ii) whether or not this can be achieved in practice, i.e., whether  $\ell^1$ -minimization produces the sparsest solution. The problem is that practically-implementable recovery conditions under which  $\ell^1$ -minimization is guaranteed to find the sparsest solution require that the vectors in the dictionary be *incoherent*, or in some way “spread out” in space. These guarantees hold with high probability, for example, on dictionaries of vectors that are randomly-generated from certain probability distributions and dictionaries consisting of randomly-selected rows of the discrete Fourier transform matrix [12, 14, 29]. Obviously, unlike these examples, data samples in the same class are often *highly-correlated*. In fact, strong inner-class similarity generally makes the data *easier* to classify.

---

We show that the assumptions under which Wright et al. argue for the effectiveness of sparsity in SRC are in direct contradiction with these recovery guarantees, leading us to question the role of sparsity in the SRC (and more broadly, classification) framework. However, using a randomly-generated database designed to model facial images, we show that  $\ell^1$ -minimization can still recover the sparsest solution on highly-correlated data, provided that the classes are sufficiently well-separated. Thus the lack of implementable equivalence *guarantee* should not automatically imply lack of *equivalence* in SRC, at least on certain databases.

Lastly, we investigate the feasibility and implementation of a transform that maximally spreads out the training samples in each class while still maintaining the data set's class structure. Though there are strict limitations on the design of such a transform, which we describe in detail in Chapter 12, we demonstrate that the higher-dimensional space can allow for the application of equivalence guarantees without negatively affecting the data set's ability to be classified. This makes it possible to examine the relationship between classification accuracy and the sparsity of the coefficient vector in SRC, and how close this is to the (provably) sparsest solution. We demonstrate that *approximate* (and not strict) equivalence between the  $\ell^1$ -minimized solution and the sparsest solution is the key to the success of SRC.

The dissertation is organized as follows: Part 1 concerns the analysis of and our proposed improvement to representation-based classification. In Chapter 2, we review the problem of classification and several key ideas necessary to understanding and assessing classification algorithms. In Chapter 3, we discuss several existing representation-based algorithms and evaluate their theoretical implications and efficacy. Chapters 4 and 5 contain the presentation and empirical evaluation of our proposed classification algorithm LPCA-SRC, a modification (or extension) of the methods presented in Chapter 3. In Chapter 6, we prove a theorem related to the success of our proposed method, and in Chapter 7, we describe and appraise some closely-related alternatives to LPCA-SRC based on the same ideology. A brief conclusion of Part 1 is contained in Chapter 8. We note that Chapters 4-7 contain original contributions.

Part 2 analyzes the use of sparsity in classification, particularly in the sparse representation framework. In Chapter 9, we motivate and review sparsity recovery guarantees, and in Chapter 10, we assess their applicability to classification data. Chapter 11 presents empirical findings relating

## 1.1. NOTATION

---

sparse recovery and class-structured data, and in Chapter 12, we investigate the feasibility of a nonlinear data transform to force the aforementioned recovery guarantees to hold. These last three chapters contain original contributions. We conclude Part 2 in Chapter 13.

Parts 1 and 2 are followed by our final remarks, and Appendix A contains the proof of the tangent bound modification discussed in Chapter 6.

### 1.1. Notation

We assume that the input data is represented by vectors in  $\mathbb{R}^m$ , and we denote the number of classes by  $L(< \infty)$ . The training set is denoted by  $X_{\text{tr}} = [\mathbf{x}_1, \dots, \mathbf{x}_{N_{\text{tr}}}] \in \mathbb{R}^{m \times N_{\text{tr}}}$ , and we refer to a given test sample by  $\mathbf{y} \in \mathbb{R}^m$ .

The  $\ell^1$ -norm of a vector  $\mathbf{x} \in \mathbb{R}^m$  with coordinates  $x_i$ ,  $1 \leq i \leq m$ , is defined as

$$\|\mathbf{x}\|_1 := \sum_{i=1}^m |x_i|,$$

and it should not be confused with the  $\ell^2$ -norm, which is given by

$$\|\mathbf{x}\|_2 := \left( \sum_{i=1}^m x_i^2 \right)^{1/2}.$$

We use the notation

$$\|\mathbf{x}\|_0 := \#\{x_i \mid x_i \neq 0\}$$

to refer to the  $\ell^0$ -“norm,” i.e., the number of nonzero coordinates of the vector  $\mathbf{x}$ . Note that this is just a pseudonorm because it does not satisfy homogeneity.

## 1.2. Lists of Figures and Tables

### Figures

3.1	An example of the per-class decomposition approach in RNS failing	27
3.2	An example of the collaborative-representation mechanism in SRC	28
3.3	Comparison of the class residuals in RNS and SRC in classification of a scarf image from the face database AR	29
5.1	A realization of the first three coordinates of the synthetic database training set	45
5.2	Box plot of the average classification accuracy of LPCA-SRC and competitive algorithms on the synthetic database with varying training class size	47
5.3	Box plot of the average classification accuracy of LPCA-SRC and competitive algorithms on the synthetic database with varying noise level	49
5.4	An example of an original image, the image recovered from PCA dimension 30, and the corresponding tangent vector in LPCA-SRC	58
5.5	An example of an original image, the image recovered from PCA dimension 30, and the corresponding tangent vector in LPCA-SRC	59
5.6	An example of an original image, the image recovered from PCA dimension 30, and the corresponding tangent vector in LPCA-SRC	59
7.1	Box plot of the average classification accuracy of the original and $\ell^2$ -regularized versions of algorithms on the synthetic database with varying training class size	80
7.2	Box plot of the average classification accuracy of the original and $\ell^2$ -regularized versions of algorithms on the synthetic database with varying noise level	81
7.3	Box plot of the average classification accuracy of the original and SSR versions of algorithms on the synthetic database with varying class size	87
7.4	Box plot of the average classification accuracy of the original and SSR versions of algorithms on the synthetic database with varying noise level	88

---

11.1	Illustration of random database model over increasing values of within-class correlation	119
11.2	Recovery results on random database model in the case of no noise	121
11.3	Asymptotic recovery results on the random database model	123
11.4	Recovery results on the random database model as the sparsity level is varied	125
11.5	Recovery results on the random database model in the case of thresholding	126
11.6	Recovery results on the random database model in the case of noise	129
12.1	Illustration of decreasing the mutual coherence of a data set by embedding it into higher dimension	134
12.2	Average sparsity, accuracy, and support quantities as correlation increased in the kernel setup	148
12.3	Median correlation between the test sample and training samples in the same class, and between the test sample and training samples in different classes	149
12.4	Average coefficient magnitudes and class contributions of coefficient vectors corresponding to class $l = 20$ test samples at various levels of coherence	150

## Tables

5.1	Mean training sample signal-to-noise ratio of the synthetic database	45
5.2	Average runtime of LPCA-SRC and compared algorithms on the synthetic database with varying training class size	48
5.3	Average accuracy and standard error of LPCA-SRC and compared algorithms on the AR face database	53
5.4	Average runtime of LPCA-SRC and compared algorithms on the AR face database	54
5.5	Average accuracy and standard error of LPCA-SRC and compared algorithms on the Extended Yale B and ORL face databases	56
5.6	Average runtime of LPCA-SRC and compared algorithms on the Extended Yale B face database	57
5.7	Average runtime of LPCA-SRC and compared algorithms on the ORL face database	58
5.8	Average energy retained in PCA dimensionality reduction to various dimensions for the face database training sets	60
7.1	Average runtime of the original and $\ell^2$ -regularized versions of algorithms on the synthetic database	82
7.2	Average accuracy, standard error, and runtime of the original and $\ell^2$ -regularized versions of algorithms on the AR face database	83
7.3	Average number of nontrivial coefficients (divided by the total number of coefficients) in the original and $\ell^2$ -regularized versions of the algorithms	83
7.4	Average runtime of the original and SSR versions of algorithms on the synthetic database	88
7.5	Average accuracy and standard error of the original and SSR versions of algorithms on the AR face database	89

---

7.6	Average number of classes (divided by the total number of classes) with nontrivial coefficients in the original and SSR versions of the algorithms	90
10.1	Average mutual coherence computed from the training samples of the face databases	108
11.1	Specification of parameters in the random database model	120
11.2	Average SRC class residuals on the random database model in the case of noise	130

## **Part 1**

# **Improving Representation-Based Classification**

## CHAPTER 2

# Classification

In this chapter, we describe the classification task and its inherent difficulties, and we define important terminology and concepts that will be used frequently in later chapters.

### 2.1. The Task

**2.1.1. What is Classification?** We, as humans, have an almost uncanny ability to recognize and classify objects around us. After making a new acquaintance, we are able to recognize her the next time we see her. We can tell the difference between an image of a building and an image of a car and identify which is which, even when only part of the object is visible or when the image has very low-resolution. For the most part, we can read each other's handwriting, even though there are countless ways to form each letter and everyone does so slightly differently.

With the rise of big data, it has become increasingly important that we find methods of passing on our natural recognition ability to machines. Having access to large amounts of data is worthless if we cannot efficiently sort it. A crime scene technician does not have time to look at millions of police photographs one-by-one in order to identify a suspect from an image pulled off a street camera. Google does not house a room full of people who manually comb through images to give us our image search results. These processes must happen automatically and incredibly quickly, and ideally, return highly-accurate results that are not subject to human error.

This task is also important because many technological advances are contingent on the automated ability to classify. A self-driven car must be able to tell the difference between pedestrians and other cars (moving or stationary) in order to avoid or minimize damage in accidents. Doctors need the technology to detect malignant tumors from MRI scans, and it is important that these are not confused with benign tumors or noise in the image. At their core, spam email filters and computer virus detectors are exactly automated classifiers, and there are countless other examples.

## 2.1. THE TASK

---

Despite its relevance, automated classification is often very difficult. Consider, for example, that for face recognition, not only must the algorithm be able to handle within-class variation in properties such as expression, face/head angle, and changes in hair or makeup, it must also be robust to differences that may occur in the image setting, most notably, the lighting conditions [71]. Further, in real-world settings, the data may suffer from noise corruption.

Let us formalize the classification problem. We say that a set of samples has *class structure* if the samples can be sorted into *classes* wherein each sample shares certain properties with all other samples in its class; further, these common properties defining the class are in some way fundamentally different than those defining other classes. In face recognition, for example, the classes correspond to the different subjects (people) in the database. For identification purposes, each class is denoted by a unique label, and every sample in the class inherits that label. Given a data set with class structure, *classification* is the task of assigning labels to samples based on a small set of (already) labeled data. We call the set of samples whose class labels are known the *training set* and any sample we want to classify a *test* sample. A *classification algorithm* is a list of programmable instructions to perform classification.

**2.1.2. Dependency on the Training Set.** In this setup, it is easy to see that the quality of the training data critically affects the performance of the classification algorithm. For example, if we give a classification algorithm several handwritten samples of the digits  $0, 1, \dots, 9$  along with their manually-assigned labels, correctness of the algorithm’s characterization of a new digit depends on the quantity and variety of these training samples. For instance, we would not expect the algorithm to classify a “2” correctly if there were only a few “2’s” in the training set, especially if those “2’s” were written very differently from our test sample. Ideally, we would supply the algorithm with a training set consisting of several samples of each digit written in a wide variety of handwriting styles.

Observe that computers are immediately at a disadvantage in the classification scenario: as humans, we have a large background of life experiences that we can draw on when making classification decisions, whereas algorithms must singularly rely on their training sets. Returning to the example of separating building and car images, an algorithm will only be able to recognize a small section of tire from a test image if there is a (portion of a) tire in a training image. In contrast,

## 2.1. THE TASK

---

humans have huge, seemingly infinite training sets to use in our classification decisions. How many cars have we seen in our lives? How many digits?

It might seem that the obvious solution is to ensure that our classification algorithm has access to a vast and diverse training set. And this is certainly the case in some applications in which training samples can be obtained cheaply. However, in many situations, it is prohibitively expensive or even impossible to obtain a large training set. Without an automated classification algorithm, the labeling of training samples must be performed by hand (i.e., by a human), and this can be extremely time-consuming. Consider, for example, cataloging hundreds of digits one-by-one. Further, in applications such as cancer detection, we might have only a handful of positive cases, i.e., MRI scans corresponding to subjects who were later treated for especially rare types of cancer. Because of these reasons, it is crucial that a widely-applicable classification algorithm be able to “make the most” of the training set, determining correct class assignments given only a limited amount of training data.

**REMARK 2.1.1.** Semi-supervised learning *is a general method of classification that uses information contained in unlabeled data (separate from the given test sample), in addition to the labeled/training data, in order to make classification decisions. This approach can be very effective in ameliorating problems resulting from the infeasibility of obtaining a large number of labeled samples [119], especially since unlabeled data is often cheaper and/or more available than labeled data. However, semi-supervised learning has its own disadvantages, e.g., the addition of a large amount of unlabeled data (generally required for this approach to be effective) may significantly increase the computational cost of performing the actual classification. Though we do not consider the semi-supervised learning approach in this dissertation, instead adhering to what is known as purely supervised learning, much of our work could potentially be modified for application in the semi-supervised framework. See, for example, the use of sparse-representation in semi-supervised learning via Cheng et al.’s  $\ell^1$ -graph [21].*

**2.1.3. Assessing Classification Algorithms.** In assessing the goodness of a classification algorithm, we consider several things. Clearly, we want the algorithm to provide us with accurate results, i.e., we want a high percentage of its classification decisions to be correct. In many

## 2.1. THE TASK

---

applications, we often need the algorithm to maintain high accuracy in the presence of significant within-class variation and noise or other data corruption. Secondly, it must be computationally efficient, so that it can be applied to large (either in number of samples or sample dimension) data sets. A probably less obvious criterion is the *simplicity* of the algorithm, in both its implementation and comprehensibility. In order for a classification algorithm to be broadly adopted, it must be able to be understood by those outside the author’s field. Lastly, we consider the *interpretability* of the classification results: does the algorithm provide us with information beyond the class assignment, such as which elements of the test sample are most important to determining its class, or the specific training samples that are most similar to it?

For illustration, consider the simplest method of classification: *k*-nearest neighbors (*k*NN) [23]. In *k*NN, the *k* closest neighbors of the given test sample (according to some distance metric) are selected from the training set, and the test sample is classified according to the most-represented class among these neighbors. When  $k = 1$ , the test sample is assigned to the class of its nearest neighbor in the training set. *k*NN can have surprisingly high classification accuracy when the distance metric used accurately captures class structure or when the data is properly pre-processed (see Section 2.2); however, it is generally quite sensitive to noise. The computational efficiency of *k*NN depends on the efficiency of computing the distance metric; when Euclidean distance is used, for example, it is relatively fast. Further, as we mentioned, it cannot be beat in terms of simplicity. However, *k*NN gives us no additional information beyond the test sample’s class assignment.

More sophisticated classification methods include naive Bayes classification, artificial neural networks, linear discriminant analysis [42, 74], and support vector machines [22], each having their own strengths and weaknesses. Due to the difficulty of the classification task, the development of algorithms that rate highly against all the above criteria is still very much an open problem.

As mentioned in the introduction, in this dissertation we focus on the *representation-based approach to classification*, which is described in detail in Chapter 3. Though this type of algorithm can be applied to digit/object classification and many of the other examples discussed in Section 2.1.1 above (recall, for instance, the various extensions of SRC mentioned in Chapter 1), it is especially well-suited for *face recognition*, and so this application will play a recurring role throughout our work. In particular, we propose a new representation-based classification algorithm in Chapter

## 2.2. PRE-PROCESSING AND FEATURE EXTRACTION

---

4, and we show in Chapters 5, 7, and 8 that it performs better against some of the above criteria than existing similar methods.

### 2.2. Pre-Processing and Feature Extraction

Before data can be classified, it must be formatted in a way that can be read by the algorithm. Though this depends on the algorithm used, we will assume that the input data is given to us as vectors, as is common. For example, if we wish to classify a set of images, we will prepare the data by taking each image and concatenating the transposed rows of pixel values into one long vector. Thus if the image is originally  $p \times p$  pixels, the resulting vector will have length  $p^2$ .

There are numerous operations on the (vectorized) data that can be performed prior to classification in order to improve the performance of the algorithm. Essentially, these processes eliminate differences in the samples that are not relevant to classification. For example, if the data samples are viewed as realizations of a random variable, *sphering* or *whitening* is the process of transforming the data so that the sample covariance matrix of the transformed data is equal to the identity matrix. This tells the algorithm that the covariance between samples should not be considered when determining class. Alternatively, *centering* the data shifts each vector by the mean vector of the data set, so that the resulting data is centered around the origin; this should be used when the algorithm automatically considers the samples' positions around **0** and we want it to be ignored. Algorithms that are sensitive to differences in the norms of the data samples or the scaling of their coordinates may benefit from sample or coordinate *normalization* when these attributes are not indicative of the samples' class structure. More complicated pre-processing methods include techniques to reduce sample noise [61] and to fill in missing values [75], e.g., those lost due to sample degradation.

Another procedure that is extremely effective in improving both the accuracy and efficiency of the classification algorithm is *feature extraction*. The simplest feature extraction methods are designed to simply pick out the coordinates of the samples that are most relevant (with respect to some criterion) to determining class. This is called *feature selection*. The irrelevant coordinates

## 2.2. PRE-PROCESSING AND FEATURE EXTRACTION

---

are determined using the training data and removed, thus helping to ensure that the subsequently-applied classification algorithm considers only information important to the task at hand. Further, by making the data samples more concise (shorter), the algorithm will run more efficiently.

More sophisticated feature extraction methods involve actually computing new features from the data samples and performing classification on these feature vectors. We describe some examples:

- *Dimensionality reduction:* Methods to perform dimensionality reduction are based on the assumption that the data lies in some lower-dimensional space, i.e., that samples originally in  $\mathbb{R}^D$  can be expressed as vectors in  $\mathbb{R}^m$ ,  $m < D$ , after a transformation to a new coordinate system. In *principal component analysis* (PCA)—one of the most popular dimensionality reduction methods—the axes of this new coordinate system are the axes of maximum variance of the training data. Alternatively, some nonlinear and more discriminative dimensionality reduction algorithms are specifically aimed at preserving the local information in the transformation, as in *locally linear embedding* [78] and *Laplacian eigenmaps* [4]. In all cases, samples are transformed into the new coordinate system and classification is performed in this new space.

The idea behind dimensionality reduction is that the representation of the data set in the new space better reveals its structure. Further, since the dimension of the new space is smaller than that of the original space, analysis of the data (e.g., classification) becomes more computationally efficient.

- *Local feature extraction:* These methods are designed to pick out portions of the data samples that contain some critical information. Once again, classification is performed using the extracted features instead of the original images. In facial recognition, for example, the eyes, nose, and mouth portions of an image contain much of the person’s identity, and so a local feature extractor might aim to classify the data set based solely on this information.

To achieve good classification in general, the extracted local features must be robust to variants such as scaling, rotation, and illumination conditions (when the sample represents an image). The *scale-invariant feature transform* (SIFT) [63], for example, determines key locations from training images by extracting feature vectors from image “blobs” that

### 2.3. THE CLASS MANIFOLD ASSUMPTION

---

appear as each image is increasingly blurred. These feature vectors are designed to be invariant to various transformations. Other local feature extractors include *histograms of oriented gradients* [24] and *local binary patterns* [70].

- *Discriminant feature extraction*: Not necessarily disjoint from the above two categories, these types of feature extractors are designed specifically for the task at hand, e.g., compression or noise removal, or in our case, classification. For example, the *local discriminant basis selection algorithm* [79, 80] can be used to select the basis (i.e., transformation matrix) from a dictionary of bases that minimizes the pairwise divergence between classes in the transform space, thus improving classification. If the dictionary of bases used is such that the expansion coefficients (i.e., the extracted features) represent local aspects of the original data samples (e.g., wavelet bases), then this method can be categorized as local feature extraction [79, 80]. Further, if only the most discriminative basis elements are used, this method performs dimensionality reduction.

Though our focus in this dissertation is on classification algorithms and not feature extraction, we remark that the choice of feature extractor can be as important as that of the classifier. Simple methods such as  $k$ NN can achieve surprisingly good performance with the right pre-processing and feature extraction methods, and classification results can usually be greatly improved by putting thought and effort into choosing the best features. This should be kept in mind in Chapters 5 and 7, in which our experiments utilize only basic PCA for dimensionality reduction.

### 2.3. The Class Manifold Assumption

One approach to designing a classification algorithm begins by assuming that the data in each class lies on (or close to) a smooth, low-dimensional manifold. There is significant precedence for this assumption: In face recognition, for example, it has been argued by Chang et al. that, since people change expression continuously over time, this assumption holds on the set of facial images that vary in expression [20]. A similar argument was made by Seung and Lee based on the smoothness of facial rotation [81]. It has also been shown that the space of all images of a fixed subject under varying illumination conditions can be well-approximated by a low-dimensional linear subspace [59]; see Chapter 11 for further details. Additional discussion and references regarding the manifold structure

### 2.3. THE CLASS MANIFOLD ASSUMPTION

---

of face databases are contained in Chapters 4 and 5. In the case of handwritten digit classification, Simard et al. modeled the set of all variations of a sample as the output of a differentiable function [85] (see Section 3.1.4), and Yang et al. applied this class manifold assumption not only to the recognition of handwritten digits but also to vehicle classification [110]. There are other examples.

In the class manifold approach to classification, the training data is used to (either explicitly or implicitly) produce a model of each class manifold, and the test sample is then classified based on a method of determining best fit. Though some of the classification algorithms we discuss do not explicitly use this procedure, we assume throughout this dissertation that the data we wish to classify has this structure. Locally, these smooth, low-dimensional class manifolds look like Euclidean space, and it is sufficient for the reader to imagine, for example, the topology of a golf course, in order to understand our use of this term in subsequent chapters.

## CHAPTER 3

# Representation-Based Classification

This chapter reviews relevant methods of representation-based classification and discusses their theoretical motivations and empirical efficacy.

### 3.1. Per-Class Decomposition

In this section, we consider representation-based classifiers that compute an approximation or representation of the test sample using the training samples in each class separately.

**3.1.1. 1-Nearest Neighbors.** *k-nearest neighbors* (kNN) can be viewed as a representation-based classifier when Euclidean distance is used and  $k = 1$ . Recall that we classify  $\mathbf{y}$  in 1NN by finding its nearest sample in the training set and then assigning  $\mathbf{y}$  to the class of that training sample. This is equivalent to solving

$$(3.1) \quad \text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \left\{ \min_{1 \leq j \leq N_l} \|\mathbf{y} - \mathbf{x}_j^{(l)}\|_2 \right\},$$

where the set  $\{\mathbf{x}_1^{(l)}, \dots, \mathbf{x}_{N_l}^{(l)}\}$  contains the  $N_l$  training samples in class  $l$ ,  $l = 1, \dots, L$ . From this minimization problem, we can see that 1NN works by determining the best approximation of  $\mathbf{y}$  using a *single* training sample.

This method is extremely popular in practice due to its simplicity. Further, it has been shown that as the number of samples  $N_{\text{tr}}$  goes to infinity, the probability of 1NN computing the wrong classification assignment is bounded above by two times the Bayes probability of error (the lowest possible error rate based on the probability distribution of the samples) [23]. However, 1NN has critical limitations. Because it considers only a single training sample at a time, it does not utilize information contained in each class as a whole. This makes 1NN very sensitive to noise, outliers, and the number of training samples in each class; the class assignment can easily change given a small perturbation of the data samples or the addition of new training data. Secondly, as we

### 3.1. PER-CLASS DECOMPOSITION

---

mentioned in Chapter 2, the classification process in 1NN tells us no more than the class assignment itself: it gives us no indication of the relationship between the test sample and its assigned class (except that it contains a training sample relatively near  $\mathbf{y}$ ).

**3.1.2. Nearest Subspace Classification.** *Nearest subspace classification* (NSC) works under the assumption that the correct class of a test sample is the class whose training set can approximate it with the least  $\ell^2$ -error. The distance between the test sample and its projection onto the span of the training samples in each class is computed, with classification based on minimal distance:

$$(3.2) \quad \text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \left\{ \min_{\boldsymbol{\alpha}_l \in \mathbb{R}^{N_l}} \|\mathbf{y} - X^{(l)} \boldsymbol{\alpha}_l\|_2 \right\}.$$

Here,  $X^{(l)} := [\mathbf{x}_1^{(l)}, \dots, \mathbf{x}_{N_l}^{(l)}]$  is the matrix that contains the training samples from class  $l$ ,  $l = 1, \dots, L$ . In the case that the columns of  $X^{(l)}$  are linearly independent (which requires that the sample dimension  $m$  is greater than the number of class  $l$  samples  $N_l$ ), the solution to Eq. (3.2) is given by

$$(3.3) \quad \text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \|\mathbf{y} - \hat{\mathbf{y}}_l\|_2^2,$$

where  $\hat{\mathbf{y}}_l := X^{(l)}((X^{(l)})^\top X^{(l)})^{-1}(X^{(l)})^\top \mathbf{y}$ .

As discussed by Zhang et al. [114], NSC is sensitive to small errors in the data samples. In face recognition, for example, if two subjects look fairly similar to each other, it may be the case that

$$\min_{\boldsymbol{\alpha}_l \in \mathbb{R}^{N_l}} \|\mathbf{y} - X^{(l)} \boldsymbol{\alpha}_l\|_2 \approx \min_{\boldsymbol{\alpha}_{l'} \in \mathbb{R}^{N_{l'}}} \|\mathbf{y} - X^{(l')} \boldsymbol{\alpha}_{l'}\|_2,$$

causing NSC to have trouble differentiating the  $l$ th and  $(l')$ th classes. Further, if the training class sizes are large enough ( $N_l > m$ ), then there is an infinite number of solutions  $\boldsymbol{\alpha}_l$  to Eq. (3.2), i.e.,  $\hat{\mathbf{y}}_l = \mathbf{y}$  in Eq. (3.3). In both cases, the performance of NSC can be improved by adding *regularization*.

**3.1.3. Regularized Nearest Subspace.** In *regularized nearest subspace* (RNS) [115], Eq. (3.2) in NSC is replaced with

$$(3.4) \quad \text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \left\{ \min_{\boldsymbol{\alpha}_l \in \mathbb{R}^{N_l}} \|\mathbf{y} - X^{(l)} \boldsymbol{\alpha}_l\|_2^2 + \lambda \|\boldsymbol{\alpha}_l\|_p \right\},$$

### 3.1. PER-CLASS DECOMPOSITION

---

where  $p \in \{1, 2\}$  and  $\lambda > 0$  is the trade-off between approximation and regularization. Here, the regularization term  $\|\alpha_l\|_p$  forces the coefficients in the approximation to be small.

As we discuss below, the test sample  $\mathbf{y}$  may not be close to the span of the training samples in the correct class, causing the projection distances in NSC and RNS to be poor indicators of class. This is especially likely to occur if the training sets are undersampled. The following method is less sensitive to this weakness.

**3.1.4. Tangent Distance Classification.** The motivation behind *tangent distance-based classification* (TDC) [18, 85, 110] is as follows: Given a sample  $\mathbf{x}_j \in \mathbb{R}^m$  and a vector of deformation parameters  $\gamma$ ,<sup>1</sup> we consider the transformed sample  $\tilde{\mathbf{x}}_j = S(\mathbf{x}_j, \gamma)$  and the set of small deformations of  $\mathbf{x}_j$  given by  $S_{\mathbf{x}_j} := \{\tilde{\mathbf{x}}_j \mid \exists \gamma \text{ for which } \tilde{\mathbf{x}}_j = S(\mathbf{x}_j, \gamma)\}$  [85]. In the case that the sample represents an image of an object, we can think of  $S_{\mathbf{x}_j}$  as the space of various “versions” of the object, for example, the space containing all the different ways the digit “2” can be written by hand.

TDC associates each class with one or more spaces  $S_{\mathbf{x}_j}$  and then aims to classify the test sample to the class associated with the closest such space. In the basic TDC algorithm, this is done by determining a linear approximation to each  $S_{\mathbf{x}_j}$  and then measuring the distance between these tangent hyperplanes and the test sample, with classification to the class with the closest hyperplane. Simard et al. [85] construct these tangent hyperplanes from a set of predetermined transformations of selected prototypes for each class. This is a good approach when there is prior knowledge available regarding the deformation parameters, as in the case of the handwritten digits used by Simard et al. [85].

However, such prior information is generally not available. To address this issue, the local TDC algorithm of Yang et al. [110] determines the nearest neighbors of the given test sample in each class, scales them so that  $\mathbf{y}$  approximates their mean, and then uses the scaled neighbors to construct a class-specific sample covariance matrix. The authors state that if the (scaled) neighbors are i.i.d. on a sufficiently-smooth manifold and the covariance matrix has full-rank, then the eigenvectors of the covariance matrix form a basis for the tangent hyperplane on the class manifold at the mean

---

<sup>1</sup>For example, consider  $\gamma \in \mathbb{R}^3$  with first coordinate representing rotation angle and second and third coordinates representing vertical and horizontal translation, respectively. These deformation parameters are applicable to tasks such as the classification of handwritten digits.

### 3.2. SPARSE REPRESENTATION-BASED CLASSIFICATION

---

of the scaled neighbors. The Euclidean distance between this hyperplane and the test sample can then be measured.

TDC is a representation-based method in the sense that this hyperplane can be viewed as the affine span of *newly-generated* training samples, i.e., approximate tangent hyperplane basis vectors. We discuss this further in Chapter 4. Though TDC methods can be effective, the accuracy of tangent vectors computed from training data is very sensitive to the amount of noise in the data samples. Further, most methods for computing these tangent vectors require user-specified parameters. Thus cross-validation must often be performed to estimate them (increasing the algorithm's overall runtime), with improper setting of these parameters negatively affecting the algorithm's performance.

We formally state two implementations of TDC in Chapter 5.

**3.1.5. Shortcomings of Per-Class Decomposition Methods.** In the case of limited training data, per-class approximations of the test sample may be poor, even over the correct class. Thus no class may distinguish itself under this framework. Even with the addition of tangent vectors (as in TDC), class-specific training data may fail to provide a good estimation of a generic test sample. Though these per-class methods generally perform more competitively when  $m < N_l$  and regularization is added (i.e., RNS), there are important applications in which we do not have the luxury of such large training sets. For example, face recognition is generally a small class size problem, with only a handful of training images available in each class. Further, the raw number of pixels in these images can be very large. Even with dimensionality reduction or feature extraction, it is challenging (and sometimes impossible) to satisfy the inequality  $m < N_l$ ,  $1 \leq l \leq L$ , while still retaining enough information in the extracted features to distinguish the classes.

The following algorithm is an alternative to per-class decomposition methods.

## 3.2. Sparse Representation-Based Classification

*Sparse representation-based classification* (SRC) [104] plays a central role in this dissertation. This method decomposes the test sample over the entire training set, solving the optimization

### 3.2. SPARSE REPRESENTATION-BASED CLASSIFICATION

---

problem

$$(3.5) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \|\boldsymbol{\alpha}\|_1, \text{ subject to } \mathbf{y} = X_{\text{tr}} \boldsymbol{\alpha}.$$

It is assumed that the training samples have been normalized to have  $\ell^2$ -norm equal to 1, so that the representation in Eq. (3.5) will not be affected by the samples' magnitudes. The use of the  $\ell^1$ -norm in the objective function is designed to approximate the  $\ell^0$ -“norm,” i.e., to aim at finding the smallest number of training samples that can accurately represent the test sample  $\mathbf{y}$ . It is argued that the nonzero coefficients in the representation will occur primarily at training samples in the same class as  $\mathbf{y}$ , so that

$$(3.6) \quad \text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \|\mathbf{y} - X_{\text{tr}} \delta_l(\boldsymbol{\alpha}^*)\|_2$$

produces the correct class assignment. Here,  $\delta_l$  is the indicator function that acts as the identity on all coordinates corresponding to samples in class  $l$  and sets the remaining coordinates to zero. In other words,  $\mathbf{y}$  is assigned to the class whose training samples contribute the most to the sparsest representation of  $\mathbf{y}$  over the entire training set.

The reasoning behind this is the following: It is assumed that the class manifolds are linear subspaces, so that if each class's training set contains a spanning set of the corresponding subspace, the test sample can be expressed as a linear combination of training samples in its ground truth class. If the number of training samples in each class is small relative to the number of total training samples  $N_{\text{tr}}$ , this representation is naturally sparse [104].

As real-world data are often corrupted by noise, the constrained  $\ell^1$ -minimization problem in Eq. (3.5) may be replaced with its regularized version

$$(3.7) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \left\{ \frac{1}{2} \|\mathbf{y} - X_{\text{tr}} \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \right\}.$$

Similar to RNS,  $\lambda$  is the trade-off between error in the approximation and the sparsity of the coefficient vector. We summarize SRC in Algorithm 1.

There have been many modifications to the original SRC algorithm addressing various issues, such as its high computational complexity [60, 114] and its requirement that training samples be normalized [116, 118]. It has also been mentioned that for SRC to work correctly in theory,

### 3.2. SPARSE REPRESENTATION-BASED CLASSIFICATION

---



---

**Algorithm 1** Sparse Representation-Based Classification (SRC) [104]

---

**Input:** Matrix of training samples  $X_{\text{tr}} \in \mathbb{R}^{m \times N_{\text{tr}}}$ , test sample  $\mathbf{y} \in \mathbb{R}^m$ , number of classes  $L$ , and error/sparsity trade-off  $\lambda$ .

**Output:** The computed class label of  $\mathbf{y}$ : `class_label( $\mathbf{y}$ )`.

- 1: Normalize each column of  $X_{\text{tr}}$  to have  $\ell^2$ -norm equal to 1.
  - 2: Use an  $\ell^1$ -minimization algorithm to solve either the constrained problem in Eq. (3.5) or the regularized problem in Eq. (3.7).
  - 3: **for** each class  $l = 1, \dots, L$ , **do**
  - 4:   Compute the norm of the class  $l$  residual:  $\text{err}_l(\mathbf{y}) := \|\mathbf{y} - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}^*)\|_2$ . Set  $\text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \{\text{err}_l(\mathbf{y})\}$ .
  - 5: **end for**
- 

the class manifolds must be disjoint so that no two classes may each emit a representation of the given test sample (see, e.g., the work of Hui et al. [53]). But in face recognition, for example, similar looking subjects (e.g., siblings) and extreme lighting conditions are two of the many reasons class manifolds may intersect. The objection to SRC is that, given that  $\mathbf{y}$  is in the intersection of two linear subspaces, the number of nonzero coefficients may not be able to distinguish the difference between the two possible representations of  $\mathbf{y}$ , resulting in poor classification. Though this argument assumes that the constrained optimization problem in Eq. (3.5) is used, an analogous objection holds for the regularized version in Eq. (3.7) when the test sample is sufficiently close to more than one class manifold.

This argument is related to another critique of SRC that questions whether or not the concept of sparsity is necessary, or even helpful, in classification. The phrase “sparsity does not enforce locality” has been used to point out that significantly-contributing training samples in the sparse representation of  $\mathbf{y}$  may be far away from  $\mathbf{y}$  with respect to distance metrics such as Euclidean distance [101]. Many arguments claim that the ability to enforce locality in a representation (by restricting large coefficients to dictionary elements close to  $\mathbf{y}$ ) is necessary for correct classification, and thus the sparseness prior in SRC alone is insufficient [101].

Though we discuss these issues, particularly the role of sparsity in SRC, more in Chapters 4, 5, 12, and 14, we make a couple of comments here:

- (1) The effectiveness of  $\ell^1$ -minimization in SRC may not be solely based on its ability to recover in many situations the sparsest solution (see Chapter 9 for conditions that guarantee this so-called  $\ell^1/\ell^0$ -equivalence). More broadly, its efficacy may lie in its ability to promote

### 3.3. SIMILAR METHODS

---

sparsity *as well as locality* in the approximation/representation [109]. Thus a built-in criterion of closeness may break ties when multiple, class-specific sparse representations of the test sample exist. For example, when the constraint  $\mathbf{1}^\top \boldsymbol{\alpha} = 1$  is added to the optimization problem in Eq. (3.5), Yang et al. [109] showed the following: If  $\bar{\mathbf{x}}_{\boldsymbol{\alpha}}$  denotes the mean of the training samples with nonzero coefficients in the representation  $\mathbf{y} = \mathbf{X}_{\text{tr}}\boldsymbol{\alpha}$  (note that in general, these training samples will not all be in the same class), then  $\ell^1$ -minimization finds the representation of  $\mathbf{y}$  such that the  $\ell^1$ -distance between  $\mathbf{y}$  and  $\bar{\mathbf{x}}_{\boldsymbol{\alpha}}$  is minimized, i.e.,  $\|\mathbf{y} - \bar{\mathbf{x}}_{\boldsymbol{\alpha}^*}\|_1$  is minimal. See Yang et al.’s paper [109] for details, as well as our findings in Chapter 12.

- (2) Let us suppose that the assumptions in SRC hold, namely, that (i) the class manifolds are linear subspaces spanned by their training samples and that (ii) the number of samples in each training class is small, relative to the total number of training samples. Further, assume that (iii) these class subspaces do not intersect (except for at the origin). Then sparsity in the context of SRC is a legitimate method to find the representation of  $\mathbf{y}$  in terms of the training samples in its ground truth class, and thus leads to correct classification. That is, assuming that the strict assumptions (i)-(iii) hold, no specific notion of locality is needed.

### 3.3. Similar Methods

**3.3.1. Collaborative Representation.** In addition to its sparsity component, SRC is fundamentally different from the per-class methods discussed in Section 3.1 in that it represents (or approximates) a test sample using all the classes at the same time. This in itself is a viable solution to the shortcomings of per-class decomposition methods [114]. In this global representation, the classes all “collaborate” in expressing  $\mathbf{y}$ , hence the term *collaborative representation*, as coined by Zhang et al. [114]. When  $m < N_{\text{tr}}$  (which is much more practical than requiring  $m < N_l$ , especially when the number of classes is large), there is an infinite number of exact representations of  $\mathbf{y}$  in terms of the set of training samples, and we can make the problem well-posed by adding a suitable regularization term or constraint on the coefficient vector. In SRC, this is  $\ell^1$ -minimization.

### 3.3. SIMILAR METHODS

---

Collaborative representation is based on the idea that samples in the data set, independent of their classes, have common properties. For example, consider the commonalities between facial images—each is roughly the same shape with eyes, a nose and mouth, and with the majority of pixels corresponding to areas such as skin and hair that are primarily devoid of significant class information. In a collaborative representation, this communal information can be shared over all the classes. In contrast, the discriminative information contained in the test image, unique to its ground truth class, must be contributed by training samples in this same class. Thus these samples (more specifically, their coefficients) stand out in the overall representation. It follows that a class assignment can be made by comparing the coefficients that correspond to samples in each class. Zhang et al. argue that this collaborative mechanism in SRC, and not necessarily its connection to sparsity, is fundamental to its success, at least on some data sets [114].

Though we extensively investigate the connection between SRC and sparsity throughout this dissertation, the following example specifically illustrates the efficacy of the collaborative representation aspect in SRC as compared to a per-class approach discussed in Section 3.1: We took the first  $L = 20$  classes of AR-2, a version of the AR Face Database [65] discussed in Chapter 5. Images in AR-2 contain faces with differing expression and lighting conditions, as well as a small number of images of each subject wearing: (i) sunglasses, or (ii) a scarf covering the nose and mouth. These are considered natural occlusions. After randomly selecting one from the 20 classes, we chose a test sample with scarf occlusion from that class, and then purposely removed all images of that subject with scarf occlusion from the training set. The idea is to try to trick the classification algorithm: considering that the scarf occlusion is quite major, will the algorithm select an incorrect class that contains a scarf image, or will it be able to identify the actual subject based on the portion of the test image that is not covered by the scarf?

We compared RNS (with  $p = 1$  in Eq. (3.4)) and SRC (with the optimization problem formulated as in Eq. (3.7)). The parameter  $\lambda$  was set to 0.001 in each method. Since both algorithms utilize the sparsity-promoting  $\ell^1$ -norm, this allowed us to directly compare per-class and collaborative approaches to representation-based classification.

RNS was unable to correctly classify the test image, and instead selected an incorrect class with a scarf-occluded image and placed a large coefficient at this image, as illustrated in Figure

### 3.3. SIMILAR METHODS

---

3.1. Because many of the classes contained scarf images in their training sets, the class residuals in RNS, given by

$$(3.8) \quad \min_{\alpha_l \in \mathbb{R}^{N_l}} \left\{ \|y - X^{(l)} \alpha_l\|_2^2 + \lambda \|\alpha_l\|_1 \right\}, \quad 1 \leq l \leq L,$$

were quite similar to each other. This indicates a low level of confidence in the classification decision, since RNS classifies  $y$  to the class with the minimum residual (see Eq. (3.4)). We illustrate this in Figure 3.3a.

In contrast, the collaborative representation mechanism in SRC was able to identify the correct class. The coefficient vector in SRC had nonzeros at a handful of scarf images over a number of classes (i.e., the classes *collaborated* to represent this portion of the test image), as well as nontrivial coefficients at multiple un-occluded images from the correct class. This is illustrated in Figure 3.2. Since the coefficients corresponding to scarf images were spread out among the classes, this latter type of coefficient—those at un-occluded images in the correct class—produced a very confident classification decision: the class with the smallest residual  $\text{err}_l(y)$  (defined in Algorithm 1) was clearly differentiated from the remaining classes, as illustrated in Figure 3.3b.

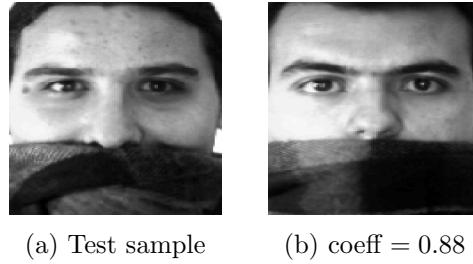


FIGURE 3.1. An example of the per-class decomposition approach in RNS failing. RNS places a very large coefficient on a scarf image from the wrong class, resulting in incorrect classification. For reference, note that all other coefficients in this class approximation had magnitudes less than 0.1.

The basic model of Zhang et al.'s *collaborative representation-based classification* (CRC) [114], designed to generalize the multi-class representation approach in SRC, is the following: Assuming that the training samples have been normalized to have  $\ell^2$ -norm equal to 1, the optimization

### 3.3. SIMILAR METHODS

---

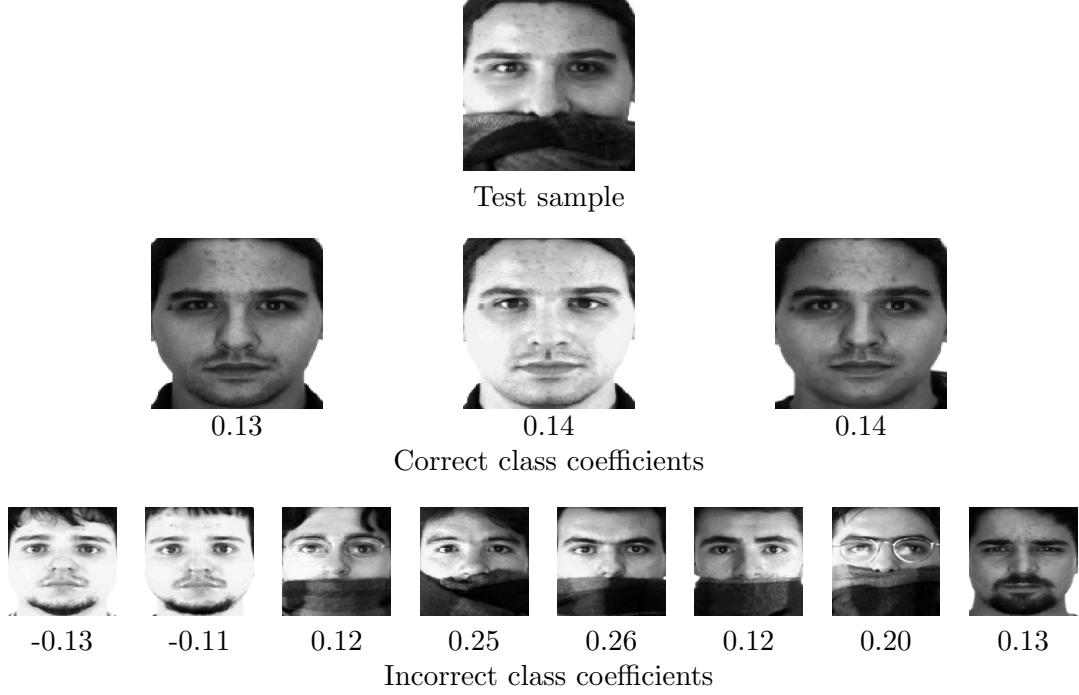


FIGURE 3.2. An example of the collaborative-representation mechanism in SRC. To decompose the occluded test sample (top row), training samples from incorrect classes collaborated to represent the scarf portion of the image (bottom row). The correct class was needed to provide class-specific details found in the un-occluded portion of the test image (middle row), resulting in correct classification. Training samples with coefficient magnitude less than 0.1 have been omitted.

problem in CRC is

$$(3.9) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \left\{ \|\mathbf{y} - X_{\text{tr}}\boldsymbol{\alpha}\|_q^2 + \lambda \|\boldsymbol{\alpha}\|_p^2 \right\},$$

with classification by

$$(3.10) \quad \text{class label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \left\{ \frac{\|\mathbf{y} - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}^*)\|_2}{\|\delta_l(\boldsymbol{\alpha}^*)\|_p} \right\}.$$

Here,  $p \in \{1, 2\}$  and  $q = 2$ , unless the data is occluded/corrupted, in which case  $q = 1$ . As in SRC,  $\delta_l$  is the indicator function that acts as the identity on all coordinates corresponding to samples in class  $l$  and sets the remaining coordinates to zero. The normalization factor  $\|\delta_l(\boldsymbol{\alpha}^*)\|_p$  in CRC is said to contain discriminative class information in addition to that contained in the class residual  $\|\mathbf{y} - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}^*)\|_2$  [114].

### 3.3. SIMILAR METHODS

---

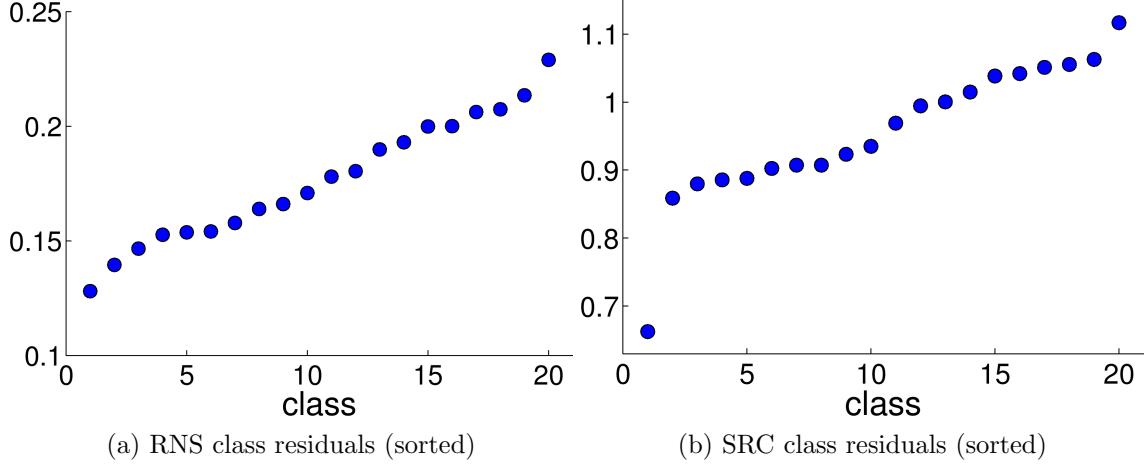


FIGURE 3.3. Comparison of the class residuals in RNS (from Eq. (3.8)) and SRC ( $\text{err}_l(\mathbf{y})$ ,  $1 \leq l \leq L$ , from Algorithm 1) in classification of a scarf image from AR-2 when no same-class scarf images are included in the training set. (a) The class residuals in RNS are quite similar to each other, indicating low confidence in the (incorrect) classification decision. (b) SRC discriminates the correct class with confidence. Classes have been reordered for best visualization.

With the exception of the normalization factor in Eq. (3.10), SRC is essentially CRC with  $p = 1$  and  $q = 2$ .<sup>2</sup> Since there is no closed-form solution to Eq. (3.9) for  $p = 1$ , it can be expensive to solve, especially for large  $m$  and  $N_{\text{tr}}$ .

Designed to address the high computational complexity of SRC, *collaborative representation-based classification with regularized least squares* (CRC-RLS) [114] is the case that  $p = q = 2$ . Though the use of the  $\ell^2$ -norm does not enforce sparsity, it serves to prevent the coefficients from blowing up. Also, as in RNS (and SRC), this regularization makes the solution more stable to small errors in the data samples.

There is a closed-form solution for  $\boldsymbol{\alpha}^*$  in CRC-RLS, given by

$$\boldsymbol{\alpha}^* = (X_{\text{tr}}^T X_{\text{tr}} + \lambda I)^{-1} X_{\text{tr}}^T \mathbf{y}.$$

In practical applications, this makes the algorithm extremely efficient. Instead of solving Eq. (3.9) for each test sample at a time, we can set  $P := (X_{\text{tr}}^T X_{\text{tr}} + \lambda I)^{-1} X_{\text{tr}}^T$  in the offline stage of CRC-RLS. Then, to classify a given test sample, we need only compute  $\boldsymbol{\alpha}^* = P \mathbf{y}$  and plug it into

<sup>2</sup>When  $p = q = 1$ , this is equivalent to what we call the *occlusion version* of SRC [104]. We describe this modification in Section 5.4.3.

### 3.3. SIMILAR METHODS

---

Eq. (3.10). Further, the CRC authors show experimentally that CRC-RLS can achieve classification rates comparable to those in the  $p = 1$  case, provided that the feature dimension  $m$  is large enough [114, 115]. In fact, in the example discussed above, CRC-RLS (like SRC) correctly classified the AR-2 test sample with scarf occlusion, due to its collaborative representation approach. Note that CRC-RLS is not specifically aimed at increasing the classification accuracy of SRC; instead, its purpose is to illustrate that there are situations in which the computational complexity of SRC can be significantly decreased without a worsening in classification decisions.

**3.3.2. Locality-Sensitive Dictionary Learning SRC.** The objective of *locality-sensitive dictionary learning for SRC* (LSDL-SRC) [101] is to modify SRC in a way that helps ensure that large coefficients in the approximation of a given test sample will be placed at training samples that are close (in terms of either  $\ell^2$ -distance or an exponential distance function) to the test sample, hence incorporating a locality constraint into the SRC algorithm.

In the dictionary learning phase of LSDL-SRC, a subset of the training samples are selected based on their ability to accurately and locally approximate a generic test sample (as implicitly approximated by the training data, so that the dictionary can be learned offline). Let  $X^{(l)} \in \mathbb{R}^{m \times N_l}$  contain the class  $l$  training samples and  $D^{(l)} \in \mathbb{R}^{m \times K_l}$  be the portion of the learned dictionary corresponding to class  $l$ . Here,  $K_l \in \mathbb{N}$ ,  $1 \leq l \leq L$ , are user-specified parameters that determine the dictionary size  $K := \sum_{l=1}^L K_l$ . Set  $\alpha_j^{(l)} \in \mathbb{R}^{K_l}$  to contain the decomposition coefficients that result when the  $j$ th training sample in the  $l$ th class is decomposed over  $D^{(l)}$ , and store these coefficient vectors in the matrix  $A^{(l)} := [\alpha_1^{(l)}, \dots, \alpha_{N_l}^{(l)}] \in \mathbb{R}^{K_l \times N_l}$ . Lastly, let  $p_j^{(l)} \in \mathbb{R}^{K_l}$  be the vector whose  $k$ th entry contains the distance between the  $j$ th training sample in the  $l$ th class and  $d_k^{(l)}$ , the  $k$ th column of  $D^{(l)}$ . With “ $\odot$ ” denoting coordinate-wise multiplication, the dictionary learning phase in LSDL-SRC [101] iteratively solves

$$(3.11) \quad \min_{D^{(l)}, A^{(l)}} \|X^{(l)} - D^{(l)}A^{(l)}\|_F^2 + \lambda_{DL} \sum_{j=1}^{N_l} \|p_j^{(l)} \odot \alpha_j^{(l)}\|_2^2 \text{ subject to } \mathbf{1}^\top \alpha_j^{(l)} = 1, \quad 1 \leq j \leq N_l,$$

for each class  $1 \leq l \leq L$ .

### 3.3. SIMILAR METHODS

---

To make this more transparent, consider that the  $k$ th coordinate of the  $j$ th term of the summation in Eq. (3.11) is

$$(3.12) \quad \text{dist}(\mathbf{x}_j^{(l)}, \mathbf{d}_k^{(l)})\alpha_{jk}^{(l)},$$

for a distance function  $\text{dist}(\cdot, \cdot)$ . The authors propose both Euclidean (or  $\ell^2$ ) distance

$$\text{dist}_{\ell^2}(\mathbf{x}_j^{(l)}, \mathbf{d}_k^{(l)}) := \|\mathbf{x}_j^{(l)} - \mathbf{d}_k^{(l)}\|_2$$

and an exponential distance function given by

$$(3.13) \quad \text{dist}_{\exp}(\mathbf{x}_j^{(l)}, \mathbf{d}_k^{(l)}) := \sqrt{\exp\left(\frac{\|\mathbf{x}_j^{(l)} - \mathbf{d}_k^{(l)}\|_2^2}{\sigma}\right)}.$$

The parameter  $\sigma$  must be specified. The second factor in Eq. (3.12),  $\alpha_{jk}^{(l)}$ , is the coefficient of  $\mathbf{d}_k^{(l)}$  when  $\mathbf{x}_j^{(l)}$  is decomposed over  $D^{(l)}$ . Thus if  $\mathbf{x}_j^{(l)}$  and  $\mathbf{d}_k^{(l)}$  are far apart, this coefficient is forced to be small, and conversely.

After solving Eq. (3.11) for each class  $1 \leq l \leq L$ , the dictionaries  $D^{(l)}$  are concatenated to form one large dictionary  $D \in \mathbb{R}^{m \times K}$ . The general SRC/CRC framework can then be used to classify a test sample  $\mathbf{y}$  using this new dictionary and replacing Eq. (3.7) with

$$(3.14) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^K} \|\mathbf{y} - D\boldsymbol{\alpha}\|_2^2 + \lambda_{\text{SRC}} \|\mathbf{p} \odot \boldsymbol{\alpha}\|_2^2.$$

Here,  $\mathbf{p} \in \mathbb{R}^K$  is the vector whose  $k$ th entry is the distance between  $\mathbf{y}$  and  $\mathbf{d}_k$ , the  $k$ th column of  $D$ . Thus significant nonzero coefficients are encouraged to occur at dictionary elements close to  $\mathbf{y}$ . After finding the minimizing argument  $\boldsymbol{\alpha}^*$  of Eq. (3.14), the final classification of  $\mathbf{y}$  is carried out using Eq. (3.6), as in SRC.

The authors of LSDL-SRC argue that the sparsity-promoting  $\ell^1$ -norm regularization term is no longer needed, as *locality implies sparsity* [101]. This allows for a closed-form solution for  $\boldsymbol{\alpha}^*$  in Eq. (3.14) (aiding computational efficiency), though the dictionary learning phase still requires iterative optimization. In experiments, LSDL-SRC has been shown to achieve higher classification accuracy than SRC, often with a smaller (more discriminative) dictionary [101]. See Chapter 5 for our empirical comparison of the two methods.

### 3.4. SUMMARY

---

#### 3.4. Summary

We presented several representation-based classification algorithms and discussed their strengths and weaknesses. Per-class decomposition methods, such as NSC, RNS, and TDC, compute an approximation of the test sample using each of the  $L$  classes separately and then classify  $\mathbf{y}$  based on the best approximation. Though these methods can be effective, all the approximations may be equally poor when  $m > N_l$ ,  $1 \leq l \leq L$ . On the other hand, the case that  $m < N_l$ ,  $1 \leq l \leq L$ , requires regularization (as in RNS), and this generous training set sampling is often an impractical assumption in applications such as face recognition.

Alternatively, collaborative representation-based methods, such as SRC, CRC-RLS, and LSDL-SRC, decompose the test sample over the entire training set at once. These algorithms only require the more feasible condition, i.e., that  $m < N_{\text{tr}}$ , and are designed to identify the correct class by locating training samples that contribute class-specific properties of the test sample to the representation. In SRC,  $\ell^1$ -minimization is used to limit the number of nonzero coefficients in the representation, whereas CRC-RLS and LSDL-SRC prevent coefficients from blowing up through the use of the more computationally-efficient  $\ell^2$ -norm. Further, LSDL-SRC is designed to restrict nonzero coefficients to training samples that are close to the given test sample, adding an aspect of locality (and some sparsity) to the global representation.

As we will see in Chapters 5 and 7, there are instances in which SRC's  $\ell^1$ -regularization is much more effective in identifying training samples in the correct class than the methods in CRC-RLS and LSDL-SRC. However, even SRC may struggle when the training set is small and the class manifolds are nonlinear, so that the correct class provides a poor approximation of the test sample. This is the motivation behind the algorithm we propose in the next chapter.

## CHAPTER 4

# Local Principal Component Analysis SRC

In this chapter, we propose a classification algorithm that is an alternative to the methods discussed in Chapter 3, in particular, SRC.<sup>1</sup>

### 4.1. Motivation

SRC explicitly makes two substantial assumptions: (i) class manifolds are *linear subspaces*, and (ii) the training data in each class *spans* its corresponding linear subspace. These assumptions are often violated. For example, facial images that vary in pose and expression lie on nonlinear class manifolds [50, 78], and as we have discussed, small training set size is one of the primary challenges in face recognition and classification as a whole.

However, SRC can still achieve high classification accuracy on data sets that violate the first of these assumptions, as demonstrated by Wright et al. [104] in experiments on the AR face database [65]. AR contains expression and occlusion variations that suggest the underlying class manifolds are nonlinear, yet SRC often outperformed SVM (support vector machines) on AR for a wide variety of feature extraction methods and feature dimensions [104]. As we have discussed, since SRC decomposes the test sample over the entire training set, components of the test sample not within the span of its ground truth class’s training samples may be absorbed by training samples from other classes. A similar fail-safe occurs when the class manifolds (linear or otherwise) are sparsely-sampled, e.g., if the second assumption is violated. However, this reliance on “off-class” training samples to partially represent or approximate the test sample can lead to misclassification, especially when the class manifolds are close together (recall the discussion in Section 3.2 regarding intersecting manifolds). In the case where class manifolds are nonlinear and/or sparsely-sampled, so that it is impossible to approximate  $\mathbf{y}$  well using only the training samples in its ground truth class, this approximation could conceivably be improved if we were able to increase the sampling density

---

<sup>1</sup>A paper on this algorithm has been submitted for publication and is available at arXiv.org [100].

## 4.2. ALGORITHM STATEMENT

---

around  $\mathbf{y}$ , “fleshing out” its local neighborhood on the (correct) class manifold. Our proposed classification algorithm aims to do just this.

### 4.2. Algorithm Statement

Our proposed algorithm, *local principal component analysis sparse representation-based classification* (LPCA-SRC), is essentially SRC with a modified dictionary. This dictionary is constructed in two steps: In the offline phase of the algorithm, we generate new training samples as a means of increasing the sampling density. Instead of the linear subspace assumption in SRC, we assume that class manifolds are well-approximated by local tangent hyperplanes. To generate new training samples, we approximate these tangent hyperplanes at individual training samples using *local principal component analysis* (local PCA), and then add (randomly-scaled and shifted versions of) the basis vectors of these tangent hyperplanes to the original training set. Naturally, the shifted and scaled tangent hyperplane basis vectors (herein referred to as “tangent vectors”) inherit the labels of their corresponding training samples. The result is an amended dictionary over which a generic test sample can ideally be decomposed using samples that approximate a local patch on the correct class manifold. In the case that the class manifolds are sparsely-sampled and/or nonlinear, this allows for a more accurate approximation of  $\mathbf{y}$  using training samples (and their computed tangent vectors) from the test sample’s ground truth class. Even in the case that class manifolds are linear subspaces, this technique ideally increases the sampling density around  $\mathbf{y}$  on its (unknown) class manifold so that it may be expressed in terms of *nearby* samples.

In the online phase of LPCA-SRC, this extended training set is “pruned” relative to the given test sample, increasing computational efficiency and the locality of the resulting dictionary. Training samples (along with their tangent vectors) are eliminated from the dictionary if their Euclidean distances to the given test sample are greater than a threshold, and then classification proceeds as in SRC as the test sample is sparsely decomposed (via  $\ell^1$ -minimization) over this local dictionary.

The method in LPCA-SRC has an additional benefit: When SRC is applied to the classification of high-resolution images (e.g.,  $> O(10^4)$  pixels), some method of dimensionality reduction is generally necessary to reduce the dimension of the raw samples, due to the high computational complexity of solving the  $\ell^1$ -minimization problem. Basic dimensionality reduction methods such

### 4.3. LOCAL PRINCIPAL COMPONENT ANALYSIS

---

as PCA (see Section 2.2) may result in the loss of class-discriminating details when the feature dimension is small. In Section 5.4.6, we show that the tangent vectors computed in LPCA-SRC can contain details of the raw images that have been lost in the PCA dimensionality reduction process.

The idea of using tangent hyperplanes for pattern recognition is not new. The most relevant instance of this is tangent distance classification (TDC) [18, 85, 110], as discussed in Chapter 3. On the other hand, there has been progress in addressing the limiting linear subspace assumption in SRC. For example, Ho et al. extended sparse coding and dictionary learning to general Riemannian manifolds [52]. Several “local” modifications of SRC implicitly ameliorate this issue; in *collaborative neighbor representation-based classification* [99] and previously-discussed LSDL-SRC, for instance, coefficients of the representation are constrained by their corresponding training samples’ distances to the test sample, and so these algorithms need only assume linearity at the local level. Other modifications have aimed at enlarging the training set in SRC specifically for face recognition, for example, using virtual images that exploit the symmetry of the human face as in both the method of Xu et al. [106] and *sample pair based sparse representation classification* [113].

Our proposed algorithm does not require explicit knowledge of the class manifolds, as in Ho et al.’s work [52], and it is applicable to classification problems outside of face recognition. Further, it is designed to improve not only the locality but also the accuracy of the approximation of the test sample in terms of its ground truth class, and its sparse representation framework gives it an empirical advantage over methods based solely on tangent distance.

We formally state the offline and online portions of LPCA-SRC in Algorithms 2 and 3, respectively. Obviously, by the definition of “offline phase,” the tangent vectors need only be computed once for any number of test samples.

#### 4.3. Local Principal Component Analysis

In LPCA-SRC (in particular, Step 5 of Algorithm 2), we use the local PCA technique of Singer and Wu [87] to compute the tangent hyperplane basis  $U^{(l,i)}$ . We outline our implementation of their method in Algorithm 4. It computes a basis for the tangent hyperplane  $T_{\mathbf{x}_i}\mathcal{M}$  at a point

#### 4.4. REMARKS ON THE CHOICE OF PARAMETERS

---

**Algorithm 2** Local PCA Sparse Representation-Based Classification (LPCA-SRC): **OFFLINE PHASE**

---

**Input:**  $X_{\text{tr}} = [\mathbf{x}_1 \dots, \mathbf{x}_{N_{\text{tr}}}] \in \mathbb{R}^{m \times N_{\text{tr}}}$ ; number of classes  $L$ ; local PCA parameters  $d$  (estimate of class manifold dimension) and  $n$  (number of neighbors).

**Output:** The normalized extended dictionary  $D \in \mathbb{R}^{m \times (N_{\text{tr}}(d+1))}$ ; pruning parameter  $r$ .

- 1: Normalize the columns of  $X_{\text{tr}}$  to have  $\ell^2$ -norm equal to 1.
- 2: **for** each class  $l = 1, \dots, L$  **do**
- 3:   Let  $\mathcal{X}^{(l)}$  be the set of class  $l$  training samples contained in  $X_{\text{tr}}$ .
- 4:   **for** each class  $l$  training sample  $\mathbf{x}_i^{(l)}, i = 1, \dots, N_l$  **do**
- 5:     Use local PCA in Algorithm 4 with set of samples  $\mathcal{X}^{(l)}$ , selected sample  $\mathbf{x}_i^{(l)}$ , and parameters  $d$  and  $n$  to compute a basis  $U^{(l,i)} := [\mathbf{u}_1^{(l,i)}, \dots, \mathbf{u}_d^{(l,i)}]$  of the tangent hyperplane approximation to the  $l$ th class manifold at  $\mathbf{x}_i^{(l)}$ . Store  $U^{(l,i)}$  and  $r_i^{(l)} := \|\mathbf{x}_{i_{n+1}}^{(l)} - \mathbf{x}_i^{(l)}\|_2$ , the distance between  $\mathbf{x}_i^{(l)}$  and its  $(n+1)$ st nearest neighbor in  $\mathcal{X}^{(l)} \setminus \{\mathbf{x}_i^{(l)}\}$ .
- 6:   **end for**
- 7: **end for**
- 8: Define the pruning parameter  $r := \text{median}\{r_i^{(l)} \mid 1 \leq i \leq N_l, 1 \leq l \leq L\}$ .
- 9: Initialize the extended dictionary  $D = \emptyset$ .
- 10: **for** each class  $l = 1, \dots, L$  **do**
- 11:   **for** each class  $l$  training sample  $\mathbf{x}_i^{(l)}, i = 1, \dots, N_l$  **do**
- 12:     Set  $c := r\gamma$ ,  $\gamma \sim \text{unif}(0, 1)$ , and form  $\tilde{X}^{(l,i)} := [c\mathbf{u}_1^{(l,i)} + \mathbf{x}_i^{(l)}, \dots, c\mathbf{u}_d^{(l,i)} + \mathbf{x}_i^{(l)}, \mathbf{x}_i^{(l)}] \in \mathbb{R}^{m \times (d+1)}$ .
- 13:     Normalize the columns of  $\tilde{X}^{(l,i)}$  to have  $\ell^2$ -norm equal to 1 and add it to the extended dictionary:  $D = [D, \tilde{X}^{(l,i)}]$ .
- 14:   **end for**
- 15: **end for**

---

$\mathbf{x}_i$  on the manifold  $\mathcal{M}$ , where it is assumed that the local neighborhood of  $\mathbf{x}_i$  on  $\mathcal{M}$  can be well-approximated by a tangent hyperplane of some dimension  $d < m$ . A particular strength of Singer and Wu's method is the weighting of neighbors by their Euclidean distances to the point  $\mathbf{x}_i$ , so that closer neighbors play a more important role in the construction of the local tangent hyperplane.

#### 4.4. Remarks on the Choice of Parameters

In this section, we detail the roles of the parameters in LPCA-SRC and suggest strategies for estimating their optimal values.

**4.4.1. Estimate of Class Manifold Dimension and Number of Neighbors.** Recall that  $d$  is the estimated dimension of each class manifold and  $n$  is the number of neighbors used in local PCA. The number of samples in the smallest training class, denoted  $N_{l_{\min}}$ , limits the range of

#### 4.4. REMARKS ON THE CHOICE OF PARAMETERS

---

**Algorithm 3** Local PCA Sparse Representation-Based Classification (LPCA-SRC): **ONLINE PHASE**

---

**Input:** Test sample  $\mathbf{y} \in \mathbb{R}^m$ ; normalized extended dictionary  $D$ ; pruning parameter  $r$ ; estimate of class manifold dimension  $d$ .

**Output:** The computed class label of  $\mathbf{y}$ : `class_label( $\mathbf{y}$ )`.

- 1: Normalize  $\mathbf{y}$  to have  $\|\mathbf{y}\|_2 = 1$ .
- 2: Initialize the pruned dictionary  $D_{\mathbf{y}} = \emptyset$  and set  $N_{\mathbf{y}} = 0$  (# of columns of  $D_{\mathbf{y}}$ ).
- 3: **for** each class  $l = 1, \dots, L$  **do**
- 4:   **for** each class  $l$  training sample  $\mathbf{x}_i^{(l)}$ ,  $i = 1, \dots, N_l$  **do**
- 5:     **if**  $\|\mathbf{y} - \mathbf{x}_i^{(l)}\|_2 \leq r$  or  $\|\mathbf{y} - (-\mathbf{x}_i^{(l)})\|_2 \leq r$  **then**
- 6:       Add the portion  $\tilde{X}^{(l,i)}$  of  $D$  corresponding to  $\mathbf{x}_i^{(l)}$  and its tangent vectors to the pruned dictionary:  $D_{\mathbf{y}} = [D_{\mathbf{y}}, \tilde{X}^{(l,i)}]$ . Assign the columns of  $\tilde{X}^{(l,i)}$  class  $l$  labels. Update  $N_{\mathbf{y}} = N_{\mathbf{y}} + (d + 1)$ .
- 7:     **end if**
- 8:   **end for**
- 9: **end for**
- 10: Use an  $\ell^1$ -minimization algorithm to compute the solution to the constrained problem

$$(4.1) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\mathbf{y}}}} \left\{ \|\boldsymbol{\alpha}\|_1 \text{ s.t. } \mathbf{y} = D_{\mathbf{y}} \boldsymbol{\alpha} \right\}$$

or the regularized problem

$$(4.2) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\mathbf{y}}}} \left\{ \frac{1}{2} \|\mathbf{y} - D_{\mathbf{y}} \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \right\}.$$

- 11: **for** each class  $l = 1, \dots, L$ , **do**
  - 12:   Compute the norm of the class  $l$  residual:  $\text{err}_l(\mathbf{y}) := \|\mathbf{y} - D_{\mathbf{y}} \delta_l(\boldsymbol{\alpha}^*)\|_2$ . Set  $\text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \{\text{err}_l(\mathbf{y})\}$ .
  - 13: **end for**
- 

**Algorithm 4** Local Principal Component Analysis (Local PCA, adapted from Singer and Wu [87])

**Input:** Set of samples  $\mathcal{X}$ , selected sample  $\mathbf{x}_i \in \mathcal{X}$ , dimension of tangent hyperplane  $d$ , number of neighbors  $n$ .

**Output:** The basis  $U^{(l,i)}$  of the approximated tangent hyperplane at the point  $\mathbf{x}_i$ .

- 1: Find the  $n + 1$  nearest neighbors (with respect to Euclidean distance) of  $\mathbf{x}_i$  in  $\mathcal{X} \setminus \{\mathbf{x}_i\}$ . Store the  $n$  nearest neighbors as columns of the matrix  $X_i := [\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_n}]$  and use the  $(n + 1)$ st nearest neighbor to define  $\epsilon_{\text{PCA}} := \|\mathbf{x}_{i_{n+1}} - \mathbf{x}_i\|_2^2$ .
  - 2: Form the matrix  $\bar{X}_i$  by centering the columns of  $X_i$  around  $\mathbf{x}_i$ :  $\bar{X}_i := [\mathbf{x}_{i_1} - \mathbf{x}_i, \dots, \mathbf{x}_{i_n} - \mathbf{x}_i]$ .
  - 3: Form a diagonal weight matrix  $D_i$  based on the distance between each neighbor and  $\mathbf{x}_i$  as follows: Let  

$$D_i(j, j) = \sqrt{K \left( \frac{\|\mathbf{x}_{i_j} - \mathbf{x}_i\|_2}{\sqrt{\epsilon_{\text{PCA}}}} \right)}, \quad j = 1, \dots, n,$$
where  $K$  is the Epanechnikov kernel given by  

$$K(u) := (1 - u^2)\chi_{[0,1]}.$$
  - 4: Form the weighted matrix  $B_i := \bar{X}_i D_i$ .
  - 5: Find the first  $d$  left singular vectors of  $B_i$  using singular value decomposition. Denote these vectors by  $\mathbf{u}_1^{(i)}, \dots, \mathbf{u}_d^{(i)}$ . Set  $U^{(l,i)} := [\mathbf{u}_1^{(i)}, \dots, \mathbf{u}_d^{(i)}]$ .
-

#### 4.4. REMARKS ON THE CHOICE OF PARAMETERS

---

values for  $d$  and  $n$  that may be used. Specifically,

$$(4.3) \quad d \leq n < N_{l_{\min}} - 1.$$

This follows from the fact that each training sample must have at least  $n + 1$  neighbors in its own class, with the dimension  $d$  of the tangent hyperplane being bounded above by the number of columns  $n$  in the weighted matrix of neighbors  $B_i$ .

While we suggest that  $n$  be set using cross-validation, there are many methods for determining  $d$ . One may use the multiscale SVD algorithm of Little et al. [62] or Ceruti et al.’s DANCo (*Dimensionality from Angle and Norm Concentration* [16]). Alternatively,  $d$  may be set using cross-validation, as we do in our experiments in Chapter 5. Empirically, we find that cross-validation often selects  $d$  smaller than the (expected) true class manifold dimension, and in these cases, increasing  $d$  from the selected value (i.e., increasing the number of tangent vectors used) does not significantly increase classification accuracy. We expect that the addition of even a small number of tangent vectors (those indicating the directions of maximum variance on their local manifolds, per the local PCA algorithm) is enough to improve the approximation of the test sample in terms of its ground truth class. Additional tangent vectors are often unneeded. Since the value of  $d$  largely affects LPCA-SRC’s computational complexity, these observations suggest that when the true manifold dimension is large, it is better to underestimate it than overestimate it.

**REMARK 4.4.1.** *Certainly, the parameters  $d$  and  $n$  could vary per class, i.e.,  $d$  and  $n$  could be replaced with  $d_l$  and  $n_l$ , respectively, for  $l = 1, \dots, L$ . In face recognition, however, if each subject is photographed under similar conditions, e.g., the same set of lighting configurations, then we expect that the class manifold dimension is approximately the same for each subject. Further, without some prior knowledge of the class manifold structure, using distinct  $d$  and  $n$  for each class may unnecessarily complicate the setting of parameters in LPCA-SRC.*

**4.4.2. Pruning Parameter.** Recall that we only include a training sample  $\mathbf{x}_i^{(l)}$  and its tangent vectors in the pruned dictionary  $D_y$  if  $\mathbf{x}_i^{(l)}$  (or its negative) is in the closed Euclidean ball  $\overline{B_m(\mathbf{y}, r)} \subset \mathbb{R}^m$  with center  $\mathbf{y}$  and radius  $r$ . Thus  $r$  is a parameter that prunes the extended dictionary  $D$  to obtain  $D_y$ . A smaller dictionary is good in terms of computational complexity, as the

#### 4.4. REMARKS ON THE CHOICE OF PARAMETERS

---

$\ell^1$ -minimization algorithm will run faster. Further, we can obtain this computational speedup without (theoretically) degrading classification accuracy: If  $\pm \mathbf{x}_i^{(l)}$  is far from  $\mathbf{y}$  in terms of Euclidean distance, then it is assumed that  $\pm \mathbf{x}_i^{(l)}$  is not close to  $\mathbf{y}$  in terms of distance along the class manifold. Thus  $\mathbf{x}_i^{(l)}$  and its tangent vectors should not be needed in the  $\ell^1$ -minimized approximation of  $\mathbf{y}$ .

A deeper notion of the parameter  $r$  is to view it as a rough estimate of the local neighborhood radius of the data set. More precisely,  $r$  estimates the distance from a sample within which its class manifold can be well-approximated by a tangent hyperplane (at that sample). Given  $X_{\text{tr}}$  and  $n$ ,  $r$  is automatically computed, as described in Algorithm 2. That is, we set  $r$  to be the median distance between each training sample and its  $(n+1)$ st nearest neighbor (in the same class), where  $n$ , the number of neighbors in local PCA, is used to implicitly define the local neighborhood. It follows that  $r$  is a robust estimate of the local neighborhood radius, as learned from the training data.

This also explains our choice for the tangent vector scaling factor  $c = r\gamma$ , where  $\gamma \sim \text{unif}(0, 1)$ . Multiplying each tangent hyperplane basis vector  $\mathbf{u}_j^{(l,i)}$ ,  $1 \leq j \leq d$ , by this constant and then shifting it by its corresponding training sample  $\mathbf{x}_i^{(l)}$  helps to ensure that the resulting tangent vector, included in the dictionary  $D_{\mathbf{y}}$  if  $\pm \mathbf{x}_i^{(l)}$  is sufficiently close to  $\mathbf{y}$ , lies in the local neighborhood of  $\mathbf{x}_i^{(l)}$  on the  $l$ th class manifold.

**REMARK 4.4.2.** *If the test sample  $\mathbf{y}$  is far from the training data, defining  $r$  as in Algorithm 2 may produce  $D_{\mathbf{y}} = \emptyset$ , i.e., there may be no training samples within that distance of  $\mathbf{y}$ . Thus to prevent this degenerate case, we use a slightly modified technique for setting  $r$  in practice. After assigning the median neighborhood radius  $r_1 := \text{median}\{r_i^{(l)} \mid 1 \leq i \leq N_l, 1 \leq l \leq L\}$ , we define  $r_2$  to be the distance between the test sample  $\mathbf{y}$  and the closest training sample (up to sign). We then define the pruning parameter  $r := \max\{r_1, r_2\}$ . In the (degenerate) case that  $r = r_2$ , the dictionary consists of the closest training sample and its tangent vectors, leading to nearest neighbor classification instead of an algorithm error. However, experimental results indicate that the pruning parameter  $r$  is almost always equal to the median neighborhood radius  $r_1$ , and so we leave this technicality out of the official algorithm statement to make it easier to interpret.*

## 4.5. COMPUTATIONAL COMPLEXITY AND STORAGE REQUIREMENTS

---

**4.4.3. Using Cross-Validation to Set Multiple Parameters.** On data sets of which we have little prior knowledge, it may be necessary to use cross-validation to set multiple parameters in LPCA-SRC. Since grid search (searching through all parameter combinations in a brute-force manner) is typically expensive, we suggest that cross-validation be applied to the parameters  $n$ ,  $\lambda$ , and  $d$ , consecutively in that order as needed. During this process, we recommend holding  $\lambda$  equal to a small, positive value (e.g., 0.001) and setting  $d = 1$  until their respective values are determined.

Our reasoning for suggesting this procedure is the following: During experiments, we found that the LPCA-SRC algorithm can be quite sensitive to the setting of  $n$ , especially when there are many samples in each training class (since there are many possible values for  $n$ ). This is expected, as the setting of  $n$  affects both the accuracy of the tangent vectors and the pruning parameter  $r$ . In contrast, LPCA-SRC is empirically fairly robust to the values of  $\lambda$  and  $d$  used, and further, setting  $d = 1$  can result in quite good performance in LPCA-SRC, even when the true dimension of the class manifolds is expected to be larger.

## 4.5. Computational Complexity and Storage Requirements

**4.5.1. Computational Complexity of SRC.** When the  $\ell^1$ -minimization algorithm HOMOTOPY [33] is used, it is easy to see that the computational complexity of SRC is dominated by this step. This complexity is  $O(N_{\text{tr}}m\kappa + m^2\kappa)$ , where  $\kappa$  is the number of HOMOTOPY iterations [107]. HOMOTOPY has been shown to be relatively fast and good for use in robust face recognition [107]. In our experiments in Chapter 5, we use it in all classification methods requiring  $\ell^1$ -minimization.

**4.5.2. Computational Complexity of LPCA-SRC.** The computational complexity of the offline phase in LPCA-SRC (Algorithm 2) is

$$(4.4) \quad O\left(m \sum_{l=1}^L N_l^2 + N_{\text{tr}}mn\right),$$

whereas that of the online phase (Algorithm 3) is

$$(4.5) \quad O\left(N_{\text{tr}}m + \frac{N_y}{d} \log\left(\frac{N_y}{d}\right) + N_y m\kappa + m^2\kappa\right).$$

## 4.5. COMPUTATIONAL COMPLEXITY AND STORAGE REQUIREMENTS

---

Recall that  $N_y$  denotes the number of columns in the pruned dictionary  $D_y$ . We note that the offline cost in Eq. (4.4) is based on the linear nearest neighbor search algorithm for simplicity; in practice there are faster methods. In our experiments, we used ATRIA (*Advanced Triangle Inequality Algorithm* [67]) via the MATLAB TSTOOL functions `nn_prepare` and `nn_search` [68]. The first function prepares the set of class  $l$  training samples  $\mathcal{X}^{(l)}$  for nearest neighbor search at the onset, with the intention that subsequent runs of `nn_search` on this set are faster than simply doing a search without the preparation function. Other fast nearest neighbor search algorithms are available, for example, *k-d tree* [5]. The cost complexity estimates of these fast nearest neighbor search algorithms are somewhat complicated, and so we do not use them in Eq. (4.4). Hence, Eq. (4.4) could be viewed as the worst-case scenario.

Offline and online phases combined, the very worst-case computational complexity of LPCA-SRC is  $O(N_{\text{tr}}^4)$ , which occurs when the second-to-last term in Eq. (4.5) dominates: i.e., when (i)  $N_y \approx (d+1)N_{\text{tr}}$  (no pruning); (ii)  $m \approx N_{\text{tr}}$  (large relative sample dimension); (iii) very large class manifold dimension estimate  $d$ , so that  $d$  is relatively close to  $N_{\text{tr}}$  (note that this requires very large  $N_l$  for  $1 \leq l \leq L$  by Eq. (4.3), which implies that  $L$  has to be very small); and (iv)  $\kappa \approx m$  (many HOMOTOPY iterations). For small  $\kappa$  and  $N_l$ ,  $1 \leq l \leq L$ , and when the pruning parameter  $r$  results in small  $N_y$  relative to  $N_{\text{tr}}$ , then the computational complexity reduces to approximately  $O(N_{\text{tr}}m)$ .

**4.5.3. Storage Requirements.** The primary difference between the storage requirements for LPCA-SRC and SRC is that the offline phase of LPCA-SRC requires storing the matrix  $D \in \mathbb{R}^{m \times (d+1)N_{\text{tr}}}$ , which has a factor of  $d+1$  as many columns as the matrix of training samples  $X_{\text{tr}} \in \mathbb{R}^{m \times N_{\text{tr}}}$  stored in SRC. Hence the storage requirements of LPCA-SRC are  $(d+1)$  times the amount of storage required by SRC.

## CHAPTER 5

# Experiments with LPCA-SRC

In this chapter, we test our proposed classification algorithm, LPCA-SRC, on one synthetic database and three popular face databases. For all data sets, we used HOMOTOPY to solve the regularized versions of the  $\ell^1$ -minimization problems, i.e., Eq. (3.7) for SRC and Eq. (4.2) for LPCA-SRC, using version 2.0 of the L1 Homotopy toolbox [1].

### 5.1. Algorithms Compared

We compared LPCA-SRC to the original SRC,  $\text{SRC}_{\text{pruned}}$  (a modification of SRC which we explain shortly), two versions of TDC (our implementations are inspired by Yang et al. [110]), LSDL-SRC [101],  $k\text{NN}$ , and  $k\text{NN}$  over extended dictionary.

- $\text{SRC}_{\text{pruned}}$ : To test the efficacy of the tangent vectors in the LPCA-SRC dictionary, this modification of SRC prunes the dictionary of original training samples using the pruning parameter  $r$ , as in LPCA-SRC.  $\text{SRC}_{\text{pruned}}$  is exactly LPCA-SRC without the addition of tangent vectors.
- *Tangent distance classification* (TDC1 and TDC2): We compared LPCA-SRC to two versions of TDC to test the importance of our algorithm’s sparse representation framework. Both of our implementations begin by first finding a pruned matrix  $D_{\mathbf{y}}^{\text{TDC}}$  that is very similar to the dictionary  $D_{\mathbf{y}}$  in LPCA-SRC. In particular,  $D_{\mathbf{y}}^{\text{TDC}}$  can be found using Algorithm 2 (omitting Step 1) and Steps 2-9 in Algorithm 3. Neither the training nor test samples are  $\ell^2$ -normalized in the TDC methods; compared to the SRC algorithms, TDC1 and TDC2 are not sensitive to the energy of the samples. We emphasize that the resulting matrix  $D_{\mathbf{y}}^{\text{TDC}}$  contains training samples that are nearby  $\mathbf{y}$ , as well as their corresponding tangent vectors.

In TDC1, we then divide  $D_{\mathbf{y}}^{\text{TDC}}$  into the “subdictionaries”  $D_{\mathbf{y}}^{(l)}$ , where  $D_{\mathbf{y}}^{(l)}$  contains the portion of  $D_{\mathbf{y}}^{\text{TDC}}$  corresponding to class  $l$ . The test sample  $\mathbf{y}$  is next projected onto the

## 5.2. SETTING OF PARAMETERS

---

space spanned by the columns of  $D_{\mathbf{y}}^{(l)}$  to produce the vector  $\hat{\mathbf{y}}^{(l)}$ , and the final classification is performed using

$$\text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \|\mathbf{y} - \hat{\mathbf{y}}^{(l)}\|_2.$$

Our second implementation, TDC2, is similar. Instead of dividing  $D_{\mathbf{y}}^{\text{TDC}}$  according to class, however, we split it up according to training sample, obtaining the subdictionaries  $D_{\mathbf{y}}^{(l,i)}$ , where  $D_{\mathbf{y}}^{(l,i)}$  contains the original training sample  $\mathbf{x}_i^{(l)}$  and its tangent vectors. It follows that each subdictionary in TDC2 has  $d + 1$  columns. The given test sample  $\mathbf{y}$  is next projected onto the space spanned by the columns of  $D_{\mathbf{y}}^{(l,i)}$  to produce  $\hat{\mathbf{y}}_i^{(l)}$ , a vector on the (approximate) tangent hyperplane at  $\mathbf{x}_i^{(l)}$ . The final classification is performed using

$$\text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \left\{ \min_{1 \leq i \leq N_l} \|\mathbf{y} - \hat{\mathbf{y}}_i^{(l)}\|_2 \right\}.$$

- *Locality-sensitive dictionary learning for SRC* (LSDL-SRC): Recall from Chapter 3 that LSDL-SRC replaces the regularization term in Eq. (3.7) of SRC with a term that forces large coefficients to occur only at dictionary elements that are close to the given test sample. We use the exponential distance function given in Eq. (3.13), as this has been shown to produce better empirical results than the  $\ell^2$ -distance function [101].
- *k-nearest neighbors classification* (*kNN*): The test sample is classified to the most-represented class from among the nearest (in terms of Euclidean distance)  $k$  training samples ( $k$  is odd).
- *k-nearest neighbors classification over extended dictionary* (*kNN-Ext*): This is *kNN* over the columns of the (full) extended dictionary that includes the original training samples and their tangent vectors. Samples are not normalized at any stage.

### 5.2. Setting of Parameters

For the synthetic database, we used cross-validation at each instantiation of the training set to choose the best parameters  $n$ ,  $\lambda$ , and  $d$  in LPCA-SRC. (Though the true class manifold dimension is known on this database, we cannot always assume that this is the case.) We optimized the parameters consecutively as described in Section 4.4.3, each over its own set of discrete values. The

### 5.3. SYNTHETIC DATABASE

---

values for  $n$  and  $d$  had to satisfy the inequalities in Eq. (4.3). We used the same approach for the parameter  $\lambda$  in SRC, the parameters  $n$  and  $\lambda$  in  $\text{SRC}_{\text{pruned}}$ , and the parameters  $n$  and  $d$  in the TDC algorithms. Finally, we used a similar procedure for the multiple parameters in LSDL-SRC (including its number of dictionary elements), and we also set  $k$  in  $k\text{NN}$  and  $k\text{NN-Ext}$  using cross-validation.

Our approach for the face databases was very similar, though in order to save computational costs, we set some parameter values according to previously published works. In particular, we set  $\lambda = 0.001$  in LPCA-SRC, SRC, and  $\text{SRC}_{\text{pruned}}$ , as was used in SRC by Waqas et al. [99]. Additionally, we set most of the parameters in LSDL-SRC to the values used by its authors [101] on the same face databases, though we again used cross-validation to determine its number of dictionary elements.

### 5.3. Synthetic Database

**5.3.1. Database Description.** The following synthetic database is easily visualized, and its class manifolds are nonlinear (though well-approximated by local tangent planes) with many intersections. Thus it is ideal for empirically comparing LPCA-SRC and SRC. The class manifolds are sinusoidal waves normalized to lie on  $S^2$ , with underlying equations given by

$$\begin{aligned}x(t) &= \cos(t + \phi), \\y(t) &= \sin(t + \phi), \\z(t) &= A \sin(\omega t).\end{aligned}$$

We set  $\omega = 3$  and  $A = 0.5$ , and we varied  $\phi$  to obtain  $L$  classes. In particular, we set  $\phi = 2\pi/(3l)$  for data in class  $1 \leq l \leq L = 4$ . For each training and test set, we generated the same number  $N_0 = N_l$ ,  $l = 1, \dots, L$ , of samples in each class by (i) regularly sampling  $t \in [0, 2\pi]$  to obtain the points  $\mathbf{p}(t) = [x(t), y(t), z(t)]^\top$ ; (ii) computing the normalized points  $\mathbf{p}(t)/\|\mathbf{p}(t)\|_2$ ; (iii) appending 50 “noise dimensions” to obtain vectors in  $\mathbb{R}^{53}$ ; (iv) adding independent random noise to each coordinate of each point as drawn from the Gaussian distribution  $\mathcal{N}(0, \eta^2)$ ; and lastly (v) re-normalizing each point to obtain vectors of length  $m = 53$  lying on  $S^{m-1}$ . We performed classification on the resulting data samples. Note that the reason why we turned the original  $\mathbb{R}^3$

### 5.3. SYNTHETIC DATABASE

---

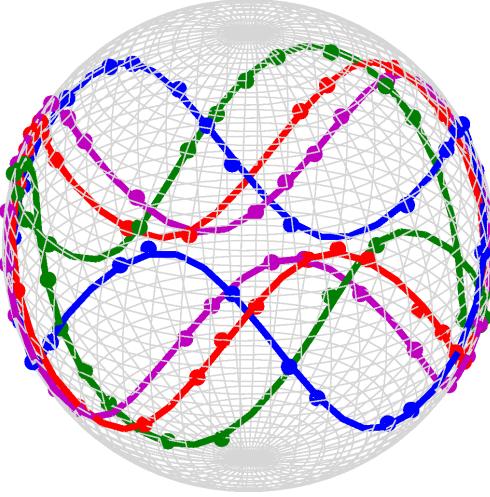


FIGURE 5.1. A realization of the first three coordinates of the synthetic database training set with  $N_0 = 25$  and  $\eta = 0.01$ . Nodes denote training samples; colors denote classes.

problem into a problem in  $\mathbb{R}^{53}$  was because SRC is designed for high-dimensional classification problems [104] and to make the problem more challenging. We emphasize that we did not apply any method of dimensionality reduction to this database.

Figure 5.1 shows the first three coordinates of a realization of the training set of the synthetic database. Note that the class manifold dimension is the same for each class and equal to 1. The signal-to-noise (SNR) ratios are displayed in Table 5.1 for  $N_0 = 25$  and various values of noise level  $\eta$ . These results were obtained by averaging the mean training sample SNR over 100 realizations of the data set.

$\eta = 0.0001$	$\eta = 0.001$	$\eta = 0.005$	$\eta = 0.01$	$\eta = 0.015$	$\eta = 0.02$	$\eta = 0.03$	$\eta = 0.05$
62.85	42.84	28.86	22.86	19.35	16.89	13.45	9.25

TABLE 5.1. Mean training sample signal-to-noise ratio (in decibels) over 100 realizations of the synthetic database with  $N_0 = 25$  and various values of noise level  $\eta$ .

**5.3.2. Experimental Results.** We performed experiments on this database, first varying the number of training samples in each class and then varying the amount of noise. The algorithms LPCA-SRC, SRC, SRC<sub>pruned</sub>, and the TDC methods significantly exceeded LSDL-SRC and the  $k$ NN methods in terms of accuracy in these experiments. In particular, these latter three methods

### 5.3. SYNTHETIC DATABASE

---

were always outperformed by LPCA-SRC by at least 10% and often by as much as 30% – 40%. Though  $k$ NN-Ext generally performed better than  $k$ NN, neither method was competitive due to its inability to distinguish individual class manifolds near intersections, a result of considering the classes in terms of a single sample (or tangent vector) at a time. On the other hand, LSDL-SRC was not *local enough*; despite its explicit locality term, this method was unable to distinguish the individual classes from within a local neighborhood of the test sample. Because of their poor performance, we do not report the results of these algorithms here.

In contrast, the approximations in LPCA-SRC, SRC, and SRC<sub>pruned</sub> often contained nonzero coefficients solely at one or two dictionary elements bordering the test sample (up to sign) on the correct class manifold. That is, these approximations were *very* sparse, and this sparsity often resulted in correct classification. The TDC methods, though generally not as competitive as these first three algorithms, also showed relatively good performance; when there was a large enough number of training samples in each class, the TDC class-specific subdictionaries were effective in discriminating between classes.

Figure 5.2 shows the average classification accuracy (over 100 trials) of the competitive algorithms as we varied the number of training samples in each class. We fixed the noise level  $\eta = 0.001$ . LPCA-SRC generally had the highest accuracy. On average, LPCA-SRC outperformed SRC by 3.5%, though this advantage slightly decreased as the sampling density increased and the tangent vectors became less useful, in the sense that there were often already enough nearby training samples in the ground truth class of  $\mathbf{y}$  to accurately approximate it without the addition of tangent vectors. SRC and SRC<sub>pruned</sub> had comparable accuracy for all tried values of  $N_0$ , indicating that the pruning parameter  $r$  was effective in removing unnecessary training samples from the SRC dictionary. Further, the increased accuracy of LPCA-SRC over SRC<sub>pruned</sub> suggests that the tangent vectors in LPCA-SRC contributed meaningful class information.

The TDC methods performed relatively poorly for small values of  $N_0$ . At low sampling densities, the TDC subdictionaries were poor models of the (local) class manifolds, leading to approximations of  $\mathbf{y}$  that were often indistinguishable from each other and resulting in poor classification. Both TDC methods improved significantly as  $N_0$  increased, with TDC2 outperforming TDC1 and in fact becoming comparable to LPCA-SRC for  $N_0 \geq 60$ . We attribute this to the extremely local nature

### 5.3. SYNTHETIC DATABASE

---

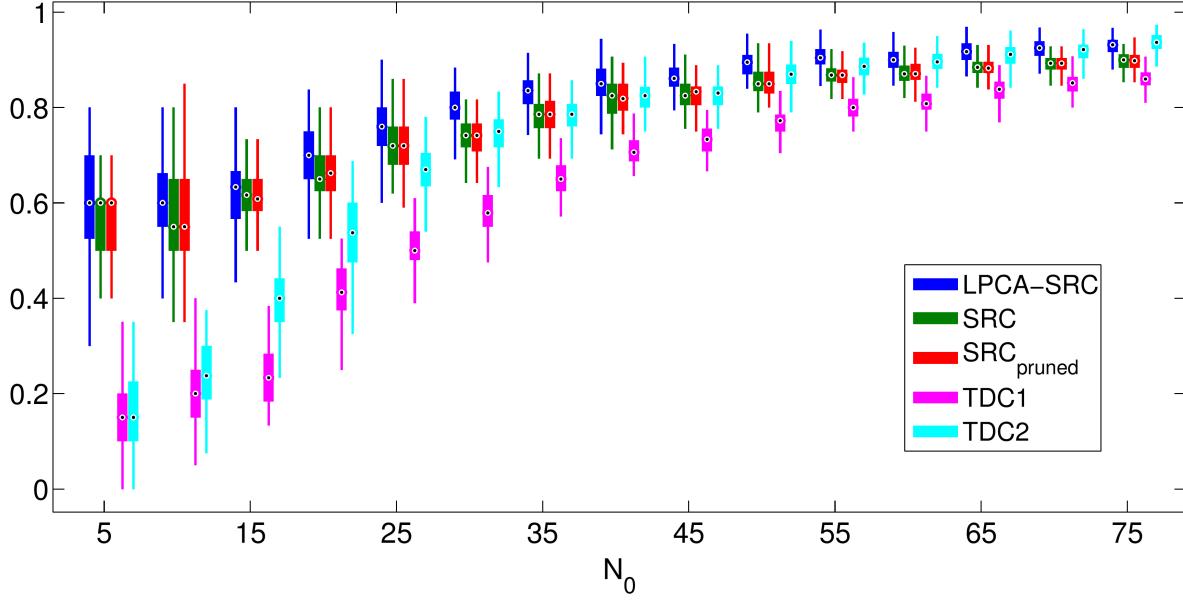


FIGURE 5.2. Box plot of the average classification accuracy (over 100 trials) of LPCA-SRC and competitive algorithms on the synthetic database with varying training class size  $N_0$ . We fixed  $\eta = 0.001$ .

of TDC2: It considers a single local patch on a class manifold at a time, rather than each class as a whole. Hence under dense sampling conditions, TDC2 effectively mimicked the successful use of sparsity in LPCA-SRC.

In Table 5.2, we display the runtime-related information of the competitive algorithms with varying training class size. In particular, we report the average runtime (in seconds), the number of columns in each algorithm’s dictionary (we refer to this as the “size” of the dictionary, as the sample dimension is fixed), and the number of HOMOTOPY iterations. These latter variables are denoted  $N$  and  $\kappa$ , respectively. The runtime does not include the time it took to perform cross-validation. For the TDC methods, we report the average subdictionary sizes, and for conciseness, we display the results for only a handful of the values of  $N_0$ . We use “N/A” to indicate that a particular statistic is not applicable to the given algorithm.

The dictionary sizes of LPCA-SRC, SRC, and  $\text{SRC}_{\text{pruned}}$  are quite informative. Recall that LPCA-SRC outperformed SRC and  $\text{SRC}_{\text{pruned}}$  (by more than 3%) for the shown values of  $N_0$ . For  $N_0 = 5$ , the dictionary in LPCA-SRC was larger than that of the two other methods, adaptively

### 5.3. SYNTHETIC DATABASE

---

Algorithm	$N_0 = 5$			$N_0 = 25$			$N_0 = 45$			$N_0 = 65$		
	t	N	$\kappa$	t	N	$\kappa$	t	N	$\kappa$	t	N	$\kappa$
LPCA-SRC	0.01	56	2	0.07	80	3	0.12	42	3	0.16	30	2
SRC	0.00	20	2	0.04	100	3	0.10	180	3	0.16	260	3
SRC <sub>pruned</sub>	0.01	20	2	0.05	79	3	0.13	146	3	0.21	201	3
TDC1	0.01	9	N/A	0.04	6	N/A	0.07	5	N/A	0.09	3	N/A
TDC2	0.02	3	N/A	0.06	2	N/A	0.09	2	N/A	0.13	2	N/A

TABLE 5.2. Average runtime in seconds (t), dictionary size ( $N$ ), and number of HO-MOTOPY iterations ( $\kappa$ ) over 100 trials of LPCA-SRC and competitive algorithms on the synthetic database with varying training class size  $N_0$ . We fixed  $\eta = 0.001$ .

retaining more samples to counter-balance the low sampling density. At large values of  $N_0$ , LPCA-SRC took full advantage of the increased sampling density, stringently pruning the set of training samples and keeping only those very close to  $\mathbf{y}$ . Due to the resulting small dictionary, it had comparable runtime to SRC despite its additional cost of computing tangent vectors. In contrast, without the addition of tangent vectors, SRC<sub>pruned</sub> was forced to keep a large number of training samples in its dictionary; the cost of the dictionary pruning step resulted in SRC<sub>pruned</sub> running slower than SRC, despite its slightly smaller dictionary. (We note that one might expect that SRC<sub>pruned</sub> would always have a smaller dictionary than LPCA-SRC since it does not include tangent vectors; this is not the case, as the value of the number-of-neighbors parameter  $n$ , and hence the pruning parameter  $r$ , may be different for the two algorithms.)

The TDC methods ran relatively fast, especially for large values of  $N_0$ . This is expected, as these algorithms do not require  $\ell^1$ -minimization.

Figure 5.3 shows the average classification accuracy (over 100 trials) of the competitive algorithms as we varied the amount of noise. We fixed  $N_0 = 25$ . (Note that we do not show the runtime results for these experiments; those for varying class size in Table 5.2 are much more revealing.) LPCA-SRC had the highest classification accuracy for low values of  $\eta$  (equivalently, when the SNR was high), outperforming SRC by as much as nearly 4%. For  $\eta \geq 0.015$  (i.e., when the SNR dropped below 20 decibels), LPCA-SRC lost its advantage over SRC and SRC<sub>pruned</sub>. This is likely due to noise degrading the accuracy of the tangent vectors. SRC and SRC<sub>pruned</sub> had nearly identical accuracy for all values of  $\eta$ ; again, this illustrates that faraway training samples (as defined by the pruning parameter  $r$ ) did not contribute to the  $\ell^1$ -minimized approximation of the test sample, and the increased accuracy of LPCA-SRC over SRC<sub>pruned</sub> for low noise values demonstrates the

### 5.3. SYNTHETIC DATABASE

---

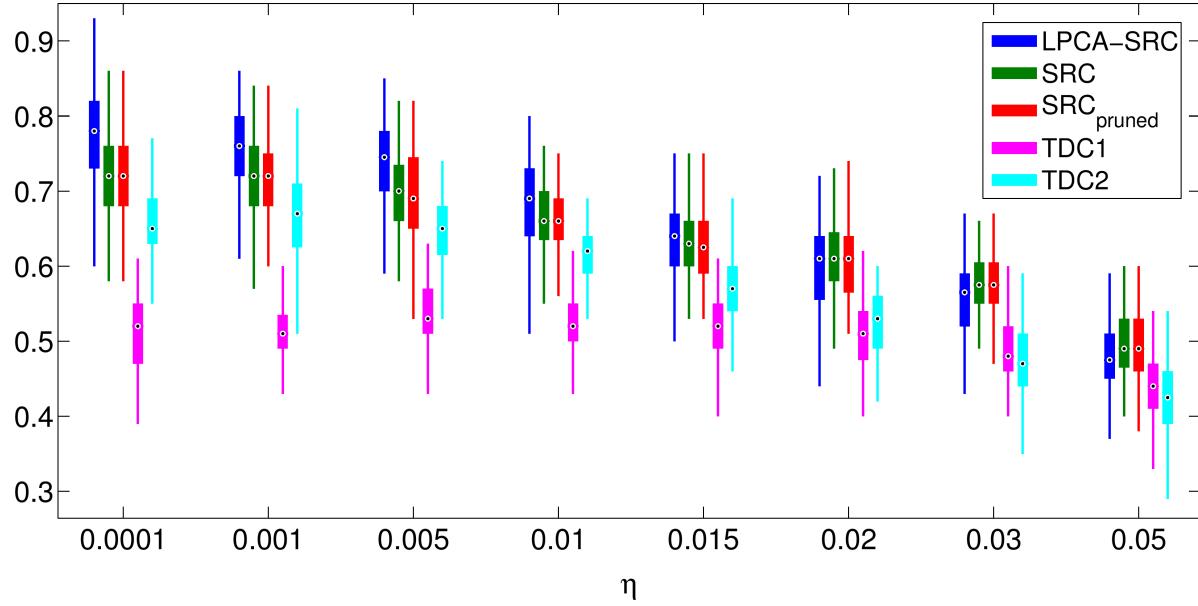


FIGURE 5.3. Box plot of the average classification accuracy (over 100 trials) of LPCA-SRC and competitive algorithms on the synthetic databases with varying noise level  $\eta$ . We fixed  $N_0 = 25$ .

efficacy of the tangent vectors in LPCA-SRC in these cases. We briefly note that when we vary the noise level for larger values of  $N_0$ , the accuracy of the tangent vectors generally improves. As a result, we see that LPCA-SRC can tolerate higher values of  $\eta$  before being outperformed by SRC and  $\text{SRC}_{\text{pruned}}$ .

TDC2 outperformed TDC1 for all but the largest values of  $\eta$ , though both algorithms were outperformed by the three SRC methods at this relatively low sampling density for the reasons discussed previously. For  $\eta \geq 0.03$ , TDC2 began performing worse than TDC1. We expect that the local patches represented by the subdictionaries in TDC2 became poor estimates of the (tangent hyperplanes of the) class manifolds as the noise increased, resulting in a decrease in classification accuracy.

In summary, the experimental results on the synthetic database show that LPCA-SRC can achieve higher classification accuracy than SRC and similar methods when the class manifolds are sparsely-sampled and the SNR is large. In these cases, the tangent vectors in LPCA-SRC help to “fill out” portions of the class manifolds that lack training samples. When the sampling density was sufficiently high, however, we saw that the tangent vectors in LPCA-SRC were less needed

## 5.4. FACE DATABASES

---

to provide an accurate, local approximation of the test sample, and thus LPCA-SRC offered a smaller advantage over SRC and  $\text{SRC}_{\text{pruned}}$ . Additionally, for higher noise (i.e., low SNR) cases, the computed tangent vectors were less reliable and the classification performance consequently deteriorated. With regard to runtime, LPCA-SRC appeared to adapt to the sampling density of the synthetic database, and though the addition of tangent vectors initially increased the dictionary size in LPCA-SRC, the online dictionary pruning step allowed for runtime comparable to SRC when the class sizes were large.

### 5.4. Face Databases

**5.4.1. Database Description.** The *AR Face Database* [65] contains 70 male and 56 female subjects photographed in two separate sessions held on different days. Each session produced 13 images of each subject, the first seven with varying lighting conditions and expressions, and the remaining six images occluded by either sunglasses or scarves under varying lighting conditions. Images were cropped to  $165 \times 120$  pixels and converted to grayscale. In our experiments, we selected the first 50 male subjects and first 50 female subjects, as was done in several papers (e.g., Wright et al. [104]), for a total of 100 classes. We performed classification on two versions of this database. The first, which we call “AR-1,” contains the 1400 un-occluded images from both sessions. The second version, “AR-2,” consists of the images in AR-1 as well as the 600 occluded images (sunglasses and scarves) from Session 1.

The *Extended Yale Face Database B* [44] contains 38 classes (subjects) with about 64 images per class. The subjects were photographed from the front under various lighting conditions. We used the version of Extended Yale B that contains manually-aligned, cropped, and resized images of dimension  $192 \times 168$ .

The *Database of Faces* (formerly “The ORL Database of Faces”) [2] contains 40 classes (subjects) with 10 images per class. The subjects were photographed from the front against dark, homogeneous backgrounds. The sets of images of some subjects contain varying lighting conditions, expressions, and facial details. Each image in ORL is initially of  $92 \times 112$  pixels.

Given existing work on the manifold structure of face databases (e.g., that of Saul and Roweis [78], He et al. [50], and Lee et al. [59]), we make the following suppositions: Since images in each

## 5.4. FACE DATABASES

---

class in AR-1 and AR-2 have extreme variations in lighting conditions and differing expressions, the class manifolds of these databases may be nonlinear. Further, the natural occlusions contained in AR-2 make these class manifolds *highly* nonlinear. Alternatively, since the images in each class in Extended Yale B differ primarily in lighting conditions, the class manifolds may be nearly linear. Lastly, since the images in some classes in ORL differ in both lighting conditions and expression, these class manifolds may be nonlinear; however, since the variations are small, these manifolds may be well-approximated by linear subspaces.

With regard to sampling density, we reiterate that Extended Yale B has large class sizes compared to AR and ORL. In our experiments, we randomly selected the same number of samples in each class to use for training, i.e., we set  $N_0 \equiv N_l$ ,  $1 \leq l \leq L$ , where  $N_0$  was half the number of samples in each class.<sup>1</sup> We used the remaining samples for testing.

**5.4.2. Dimensionality Reduction.** To perform dimensionality reduction on the face databases, we used (global) PCA to transform the raw images to  $m_{\text{PCA}} \in \{30, 56, 120\}$  dimensions before performing classification. Similar values for  $m_{\text{PCA}}$  were used by Wright et al. [104]. For the remainder of this chapter as well as Chapter 7, we will refer to the PCA-compressed versions of the raw face images as “feature vectors” and  $m_{\text{PCA}}$  as the “feature dimension.” We note that the data was not centered (around the origin) in the PCA transform space.

**5.4.3. Handling Occlusion.** Since AR-2 contains images with occlusion, we considered using the “occlusion version” of SRC (with analogous modifications to LPCA-SRC and  $\text{SRC}_{\text{pruned}}$ ) on this database. As discussed by Wright et al. [104], this model assumes that  $\mathbf{y}$  is the summation of the (unknown) true test sample  $\mathbf{y}_0$  and an (unknown) sparse error vector. The resulting modified  $\ell^1$ -minimization problem consists of appending the dictionary of training samples with the identity matrix  $I \in \mathbb{R}^{m \times m}$  and decomposing  $\mathbf{y}$  over this augmented dictionary. For more details, see Section 3.2 of the SRC paper [104]. It can be shown that the occlusion version of SRC is equivalent to CRC (see Section 3.3.1) with  $p = q = 1$ .

However, the context in which Wright et al. use the occlusion version of SRC on the AR database is critically different than our experimental setup here [104]. In the SRC paper, the samples with occlusion make up the test set. In our case, both the training and test set contain samples with and

---

<sup>1</sup>Since the class sizes vary slightly in Extended Yale B, we set  $N_0 = 32$  on this database.

## 5.4. FACE DATABASES

---

without occlusion. As a consequence, occluded samples in the training set can be used to express test samples with occlusion, and on the other hand, the use of the identity matrix to extend the dictionary in SRC results in too much error allowed in the approximation of un-occluded samples. Correspondingly, we see much worse classification performance in SRC when we use its occlusion version on AR-2. Hence, we stick to Algorithm 1 (the original version of SRC) on all face databases.

**5.4.4. AR Face Database Results.** Table 5.3 displays the average accuracy and standard error over 10 trials for the two versions of AR. LPCA-SRC had substantially higher classification accuracy than the other methods on both versions of AR with  $m_{\text{PCA}} = 30$ . This suggests that the tangent vectors in LPCA-SRC were able to recover important class information lost in the stringent PCA dimensionality reduction. As  $m_{\text{PCA}}$  increased, however, the methods SRC, SRC<sub>pruned</sub>, and LSDL-SRC became more competitive, as more discriminative information was retained in the feature vectors and less needed to be provided by the LPCA-SRC tangent vectors. SRC<sub>pruned</sub> had comparable accuracy to SRC, indicating that, once again, training samples could be removed from the SRC dictionary using the pruning parameter  $r$  without decreasing classification accuracy. In some cases, the removal of these faraway training samples slightly improved class discrimination.

For the most part, the other algorithms performed poorly on AR. The exception was LSDL-SRC, which had comparable accuracy to LPCA-SRC for  $m_{\text{PCA}} = 120$  (slightly outperforming it for AR-1) and beat SRC on AR-1 for  $m_{\text{PCA}} = 56$ . However, LSDL-SRC had lower accuracy than the SRC algorithms for  $m_{\text{PCA}} = 30$  on both versions of this database. In contrast, the TDC methods performed relatively better for  $m_{\text{PCA}} = 30$  than for larger values of  $m_{\text{PCA}}$  due to their more effective use of tangent vectors at this small feature dimension. Overall, however, their class-specific dictionaries were not as effective on this nonlinear, sparsely-sampled database as the multi-class dictionaries of the previously-discussed algorithms. Further, TDC2 often had notably high standard error, presumably because of its sensitivity to the value of the manifold dimension estimate  $d$ . This could perhaps be mitigated by using a different cross-validation procedure. Lastly,  $k\text{NN}$  and  $k\text{NN-Ext}$  had the lowest classification accuracies, though  $k\text{NN-Ext}$  offered a slight improvement over  $k\text{NN}$ . Both methods consistently selected  $k = 1$  during cross-validation.

## 5.4. FACE DATABASES

---

Algorithm	AR-1						AR-2					
	$m_{\text{PCA}} = 30$		$m_{\text{PCA}} = 56$		$m_{\text{PCA}} = 120$		$m_{\text{PCA}} = 30$		$m_{\text{PCA}} = 56$		$m_{\text{PCA}} = 120$	
	Acc	SE	Acc	SE	Acc	SE	Acc	SE	Acc	SE	Acc	SE
LPCA-SRC	<b>0.8663</b>	4.1	<b>0.9544</b>	2.3	0.9711	1.7	<b>0.7328</b>	6.0	<b>0.8844</b>	3.6	<b>0.9512</b>	2.6
SRC	0.8273	4.2	0.9357	2.6	0.9631	1.6	0.6945	4.3	0.8713	2.0	0.9450	2.4
SRC <sub>pruned</sub>	0.8277	4.8	0.9353	3.8	0.9651	1.8	0.7092	4.0	0.8781	2.7	0.9459	2.5
TDC1	0.8046	6.5	0.8430	5.1	0.8634	5.9	0.6899	4.1	0.7603	3.7	0.7985	4.5
TDC2	0.7549	19.4	0.8137	11.9	0.8303	15.4	0.6422	16.1	0.7386	3.3	0.7735	4.4
LSDL-SRC	0.8184	4.0	0.9424	2.0	<b>0.9756</b>	0.9	0.6585	5.6	0.8610	2.2	0.9498	2.9
<i>k</i> NN	0.5846	4.6	0.6301	8.0	0.6461	4.9	0.4100	3.0	0.4297	5.0	0.4554	3.2
<i>k</i> NN-Ext	0.6036	4.5	0.6487	8.2	0.6677	4.7	0.4311	3.7	0.4526	2.9	0.4794	5.7

TABLE 5.3. Average accuracy and standard error ( $\times 10^{-3}$ ) of LPCA-SRC and compared algorithms over 10 trials on AR.

Table 5.4 displays the average runtime and related results (over 10 trials) of the various classification algorithms for both versions of AR. The “dictionary size”  $N$  for  $k$ NN and  $k$ NN-Ext refers to the average size of the set from which the  $k$ -nearest neighbors are selected (e.g., for  $k$ NN,  $N = N_{\text{tr}}$ ).

The generally large dictionary sizes of LPCA-SRC (and its consequently long runtimes) indicate that minimal dictionary pruning often occurred. Thus LPCA-SRC was generally slower than SRC and SRC<sub>pruned</sub>. However, on AR-2 with  $m_{\text{PCA}} = 30$ , LPCA-SRC was able to eliminate many training samples from its dictionary, due to its effective use of tangent vectors on the (presumably) highly-nonlinear class manifolds of AR-2. At this low feature dimension, the computed tangent vectors contained more class discriminative information than nonlocal training samples, likely allowing for a more accurate—and local—approximation of  $\mathbf{y}$  on its ground truth class manifold. LPCA-SRC was faster than SRC and SRC<sub>pruned</sub> (which kept a large number of training samples) in this case, and this is impressive, considering that LPCA-SRC also outperformed these methods by nearly 4% and more than 2%, respectively.

Despite not requiring  $\ell^1$ -minimization, the TDC methods were often the slowest algorithms on the AR databases. We suspect that this is largely due to the relatively large number of classes in AR—recall that both TDC methods must compute least squares solutions (in TDC2, sometimes many of them) for each class represented in the pruned dictionary  $D_{\mathbf{y}}^{\text{TDC}}$ . Further, TDC2 selected a relatively large value of  $d$  during cross-validation (presumably so that its subdictionaries would contain a wider “snapshot” of the class manifolds), which made it even less efficient. The runtime of LSDL-SRC, unlike those of most of the other algorithms, was fairly insensitive to the feature

## 5.4. FACE DATABASES

---

dimension, and as a result, LSDL-SRC was relatively efficient for  $m_{\text{PCA}} \in \{56, 120\}$ . However, the expense of its dictionary learning phase for  $m_{\text{PCA}} = 30$ , at which the  $\ell^1$ -minimization algorithm in the SRC methods could be solved efficiently, resulted in LSDL-SRC’s relatively slow runtime. Both  $k\text{NN}$  methods ran significantly faster than all the other methods.

Algorithm	AR-1								
	$m_{\text{PCA}} = 30$			$m_{\text{PCA}} = 56$			$m_{\text{PCA}} = 120$		
	t	N	$\kappa$	t	N	$\kappa$	t	N	$\kappa$
LPCA-SRC	7.25	435	61	12.50	676	87	19.07	795	112
SRC	6.11	700	51	8.87	700	72	13.57	700	99
$\text{SRC}_{\text{pruned}}$	3.76	231	39	5.10	226	49	6.90	232	60
TDC1	11.82	16	N/A	14.24	16	N/A	24.30	19	N/A
TDC2	8.89	5	N/A	16.79	5	N/A	36.68	5	N/A
LSDL-SRC	7.78	440	N/A	8.55	470	N/A	9.72	490	N/A
$k\text{NN}$	0.01	700	N/A	0.01	700	N/A	0.02	700	N/A
$k\text{NN-Ext}$	0.08	2170	N/A	0.09	2240	N/A	0.16	2660	N/A
AR-2									
Algorithm	$m_{\text{PCA}} = 30$			$m_{\text{PCA}} = 56$			$m_{\text{PCA}} = 120$		
	t	N	$\kappa$	t	N	$\kappa$	t	N	$\kappa$
	10.53	478	58	35.27	1593	10	56.17	1690	151
SRC	11.39	1000	58	17.67	1000	85	27.74	1000	121
$\text{SRC}_{\text{pruned}}$	11.12	788	54	16.63	775	77	24.88	767	107
TDC1	20.56	25	N/A	27.51	26	N/A	43.07	26	N/A
TDC2	20.93	6	N/A	47.57	6	N/A	103.80	6	N/A
LSDL-SRC	22.70	750	N/A	16.34	620	N/A	22.19	710	N/A
$k\text{NN}$	0.02	1000	N/A	0.02	1000	N/A	0.04	1000	N/A
$k\text{NN-Ext}$	0.13	4300	N/A	0.15	3600	N/A	0.29	4400	N/A

TABLE 5.4. Average runtime in seconds (t), dictionary size (N), and number of HOMOTOPY iterations ( $\kappa$ ) over 10 trials of LPCA-SRC and compared algorithms on AR.

### 5.4.5. Extended Yale Face Database B and Database of Faces (“ORL”) Results.

Table 5.5 displays the average accuracy and standard error for Extended Yale B (over 10 trials) and ORL (over 50 trials). On Extended Yale B, LPCA-SRC had the highest accuracy for all  $m_{\text{PCA}}$ , though as we saw on the AR database, this advantage decreased as  $m_{\text{PCA}}$  increased and SRC became more competitive. SRC and  $\text{SRC}_{\text{pruned}}$  had very similar accuracy, indicating that training samples excluded from the dictionary via the pruning parameter  $r$  did not provide class information in the SRC framework. TDC1 and TDC2 had consistently mediocre performance, neither one outperforming the other over all settings of  $m_{\text{PCA}}$ , and LSDL-SRC improved as  $m_{\text{PCA}}$  increased, analogous to its behavior on AR. However, LSDL-SRC was clearly outperformed by LPCA-SRC even for  $m_{\text{PCA}} = 120$ , suggesting that the improved approximations in LPCA-SRC via its use

## 5.4. FACE DATABASES

---

of tangent vectors were more effective (even at this high feature dimension) than the procedure in LSDL-SRC.<sup>2</sup> Along these same lines, the tangent vectors in  $k$ NN-Ext offered a considerable improvement over  $k$ NN, though once again both methods reported lower accuracy than all the other algorithms. As on AR, the  $k$ NN methods consistently selected  $k = 1$  during cross-validation.

On ORL, LPCA-SRC and SRC<sub>pruned</sub> had comparable accuracy and outperformed SRC. This indicates that: (i) the pruning parameter  $r$  in LPCA-SRC and SRC<sub>pruned</sub> was *helpful* to classification (instead of simply being benign); and (ii) the tangent vectors computed in LPCA-SRC were not. With regard to (i), it must be the case that faraway training samples—those in different classes from the test sample—contributed significantly to the approximation of the test sample in SRC, negatively affecting classification performance. This is an example of *sparsity not necessarily leading to locality* (as it is relevant to class discrimination), as discussed in the LSDL-SRC paper [101]. With regard to (ii), we suspect that the tangent vectors in LPCA-SRC were simply *unneeded* to improve the classification performance on ORL. Though the approximations in SRC contained nonzero coefficients at training samples not in the same class as  $\mathbf{y}$ —presumably because of the sparse sampling and nonlinear structure of the class manifolds—many of these wrong-class training samples could be eliminated simply based on their distance to  $\mathbf{y}$ . This suggests that ORL’s class manifolds can be fairly well-separated via Euclidean distance. An additional reason for (ii) was because the PCA transform to the dimensions specified in this experiment did not result in a loss of too much information, at least compared to AR and Extended Yale B. See Table 5.8 at the end of Section 5.4.6 for this comparison.

All of the remaining methods performed relatively well on ORL. The accuracies of TDC1 and TDC2 were similar and comparable to those of SRC. We ascertained that the success of the TDC methods was not due to their use of tangent vectors but instead the result of their “per-class” approximations of the test sample. This approach was very effective on the (presumably) well-separated class manifolds of ORL. Strikingly, the accuracy of LSDL-SRC was relatively low for  $m_{\text{PCA}} = 120$ , opposite to the trend we saw on the previous face databases. The performance of LSDL-SRC could be improved for  $m_{\text{PCA}} = 120$  on this database if the samples were centered (around the origin) after PCA dimensionality reduction. However, we confirmed that LSDL-SRC

---

<sup>2</sup>LSDL-SRC was also outperformed by SRC on Extended Yale B in our experiments, in contrast to the experiments by Wei et al. in which  $m_{\text{PCA}} = 300$  [101].

## 5.4. FACE DATABASES

---

was still outperformed by LPCA-SRC in this case (albeit by a smaller margin), and its performance with centering on the other face databases was much worse than our reported results. In contrast to the results on Extended Yale B,  $k$ NN-Ext only provided a slight increase in accuracy over  $k$ NN, with the tangent vectors mimicking their unnecessary role in LPCA-SRC on this database. The value  $k = 1$  was consistently selected by both  $k$ NN and  $k$ NN-Ext during cross-validation.

Algorithm	Extended Yale B						ORL					
	$m_{\text{PCA}} = 30$		$m_{\text{PCA}} = 56$		$m_{\text{PCA}} = 120$		$m_{\text{PCA}} = 30$		$m_{\text{PCA}} = 56$		$m_{\text{PCA}} = 120$	
	Acc	SE	Acc	SE	Acc	SE	Acc	SE	Acc	SE	Acc	SE
LPCA-SRC	<b>0.9049</b>	2.9	<b>0.9530</b>	1.7	<b>0.9710</b>	1.6	<b>0.9507</b>	24.0	<b>0.9600</b>	18.0	0.9602	17.0
SRC	0.8803	2.6	0.9371	2.8	0.9633	1.4	0.9374	24.0	0.9437	22.9	0.9422	18.8
$\text{SRC}_{\text{pruned}}$	0.8804	2.7	0.9371	2.6	0.9635	1.5	0.9506	23.7	0.9580	23.5	<b>0.9605</b>	19.3
TDC1	0.8568	10.0	0.9285	2.0	0.9446	2.8	0.9364	27.1	0.9457	25.4	0.9455	21.1
TDC2	0.8826	3.9	0.9093	2.8	0.9283	3.5	0.9351	29.8	0.9429	31.1	0.9418	23.3
LSDL-SRC	0.7495	4.8	0.8774	2.5	0.9492	2.0	0.9358	25.2	0.9515	19.9	0.9251	24.3
$k$ NN	0.4300	3.5	0.5346	2.6	0.6245	3.8	0.9332	26.3	0.9387	24.7	0.9396	23.2
$k$ NN-Ext	0.5464	6.9	0.6321	5.6	0.7058	5.4	0.9338	30.5	0.9412	28.9	0.9386	23.9

TABLE 5.5. Average accuracy and standard error ( $\times 10^{-4}$ ) of LPCA-SRC and compared algorithms on Extended Yale B (over 10 trials) and ORL (over 50 trials).

Tables 5.6 and 5.7 show the runtime and related results for the Extended Yale B and ORL experiments, respectively. LPCA-SRC had much longer runtimes than SRC on Extended Yale B, especially as  $m_{\text{PCA}}$  increased. This was due to a combination of large values for  $d$  selected during cross-validation and the tangent vectors' decreasing efficacy at larger feature dimensions. However, the dictionary pruning procedure in LPCA-SRC actually eliminated a large number of training samples for all  $m_{\text{PCA}}$ ; once again, the computed tangent vectors contained more class-discriminating information than the eliminated nonlocal training samples, especially at lower feature dimensions for which details provided by these tangent vectors were especially needed. The (presumed) linearity of the class manifolds of Extended Yale B, combined with this database's relatively dense sampling, lent itself well to the accurate computation of tangent vectors—part of the reason why LPCA-SRC used so many of them. Viewing these points as newly-generated and nearby training samples, LPCA-SRC's boost in accuracy over SRC can be viewed as an argument for locality in classification. We note that we might be able to decrease the value of  $d$  in LPCA-SRC while still maintaining an advantage over SRC (see the discussion in Section 4.4.1); our cross-validation procedure is designed to obtain the highest accuracy with no regard to computational cost.

## 5.4. FACE DATABASES

---

On Extended Yale B, the TDC methods ran relatively more quickly (compared to the other algorithms) than on AR, presumably due to the much smaller number of classes on this database; both had runtimes typically between those of LPCA-SRC and SRC. Again, we see that LSDL-SRC had a relatively slow runtime for  $m_{\text{PCA}} = 30$  and became more competitive as  $m_{\text{PCA}}$  increased. Though both  $k\text{NN}$  and  $k\text{NN-Ext}$  were very fast, the large “dictionary sizes” in  $k\text{NN-Ext}$  made this algorithm clearly the slower of the two methods.

On ORL, LPCA-SRC and SRC had comparable runtimes, a result of rigorous dictionary pruning in LPCA-SRC. This algorithm and  $\text{SRC}_{\text{pruned}}$  retained roughly the same number of training samples in their respective dictionaries, and the latter was notably fast, running in about half the time as SRC. The remaining algorithms were even more efficient. TDC1 and TDC2 had comparable runtimes, both running faster than LSDL-SRC. As before,  $k\text{NN}$  and  $k\text{NN-Ext}$  had the fastest runtimes; the former was faster than the latter.

Algorithm	$m_{\text{PCA}} = 30$			$m_{\text{PCA}} = 56$			$m_{\text{PCA}} = 120$		
	t	N	$\kappa$	t	N	$\kappa$	t	N	$\kappa$
LPCA-SRC	29.20	1922	75	72.12	3359	120	141.97	3785	182
SRC	15.58	1216	62	24.70	1216	91	41.94	1216	137
$\text{SRC}_{\text{pruned}}$	15.92	1111	61	23.81	1112	88	40.50	1115	131
TDC1	8.10	20	N/A	27.62	59	N/A	42.83	59	N/A
TDC2	11.68	6	N/A	23.51	6	N/A	56.01	6	N/A
LSDL-SRC	67.30	1186	N/A	53.03	1003	N/A	38.73	821	N/A
$k\text{NN}$	0.02	1216	N/A	0.03	1216	N/A	0.05	1216	N/A
$k\text{NN-Ext}$	0.17	5350	N/A	0.25	4742	N/A	0.44	4864	N/A

TABLE 5.6. Average runtime in seconds (t), dictionary size (N), and number of HOMOTOPY iterations ( $\kappa$ ) over 10 trials of LPCA-SRC and compared algorithms on Extended Yale B.

**5.4.6. Tangent Vectors and PCA Feature Dimension.** In this section, we offer evidence to support our claim that the tangent vectors in LPCA-SRC can recover discriminative information lost during PCA transforms to low dimensions. Thus LPCA-SRC can offer a clear advantage over SRC in these cases, as we saw in experimental results on AR and Extended Yale B.

In Figures 5.4-5.6, we display three versions of three example images from AR-1. The first version is the original image (before PCA dimensionality reduction), the second version is the recovered image from PCA dimensionality reduction to dimension  $m_{\text{PCA}} = 30$ , and the third version is the recovered corresponding tangent vector computed in LPCA-SRC. In each case, the

#### 5.4. FACE DATABASES

---

Algorithm	$m_{PCA} = 30$			$m_{PCA} = 56$			$m_{PCA} = 120$		
	t	N	$\kappa$	t	N	$\kappa$	t	N	$\kappa$
LPCA-SRC	0.54	59	26	0.73	72	34	1.22	111	50
SRC	0.85	200	40	1.34	200	57	2.09	200	81
SRC <sub>pruned</sub>	0.25	19	12	0.34	26	16	0.53	39	24
TDC1	0.12	1	N/A	0.16	1	N/A	0.34	1	N/A
TDC2	0.12	3	N/A	0.23	3	N/A	0.53	3	N/A
LSDL-SRC	1.04	116	N/A	1.09	121	N/A	0.93	102	N/A
kNN	0.01	200	N/A	0.01	200	N/A	0.01	200	N/A
kNN-Ext	0.03	568	N/A	0.03	592	N/A	0.04	568	N/A

TABLE 5.7. Average runtime in seconds (t), dictionary size (N), and number of HOMOTOPY iterations ( $\kappa$ ) over 50 trials of LPCA-SRC and compared algorithms on ORL.

tangent vector contains details of the original image not found in the recovered image, supporting our claim that the tangent vectors in LPCA-SRC can recover some (but not all) of the information lost in stringent PCA dimensionality reduction.

Towards quantifying what we mean by “stringent,” Table 5.8 lists the average energy (over 10 trials) retained in the first  $m_{PCA}$  left-singular vectors of the face database training sets, along with the percent improvement in the accuracy of LPCA-SRC with respect to that of SRC and SRC<sub>pruned</sub>. We reiterate that the addition of tangent vectors did not increase classification accuracy on ORL. Taking this into account, we see a correlation between the efficacy of tangent vectors in LPCA-SRC and the stringency of the PCA dimensionality reduction.

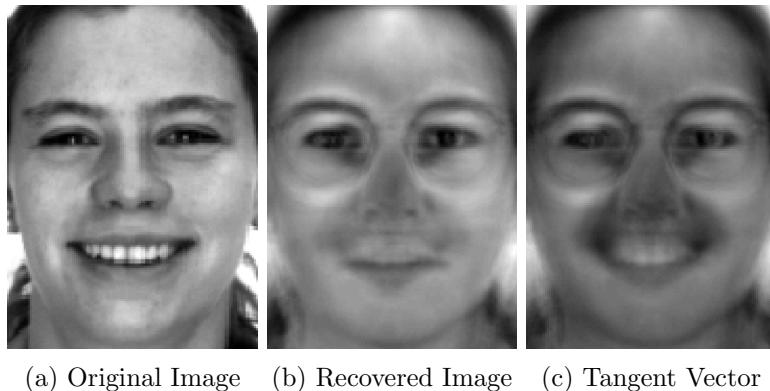


FIGURE 5.4. The tangent vector does a much better job of displaying facial details conveying “happiness” (displayed in the original image) than the recovered image. Images (b) and (c) were recovered from PCA dimension  $m_{PCA} = 30$ .

## 5.4. FACE DATABASES

---

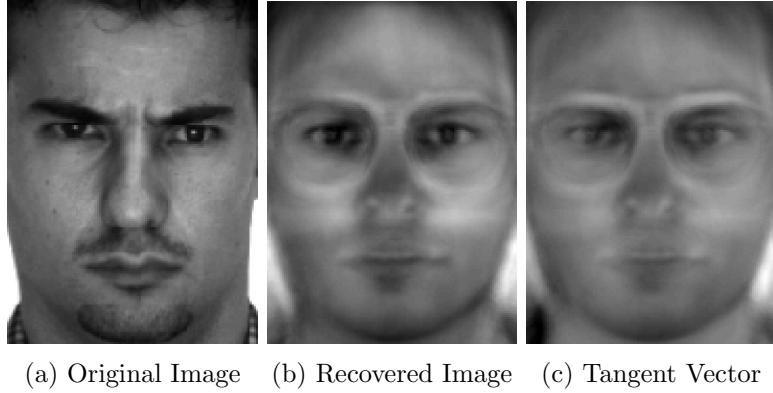


FIGURE 5.5. The tangent vector does a better job of displaying facial details conveying “anger” (displayed in the original image) than the recovered image, most notably in the subject’s eyes and eyebrows. Images (b) and (c) were recovered from PCA dimension  $m_{PCA} = 30$ .

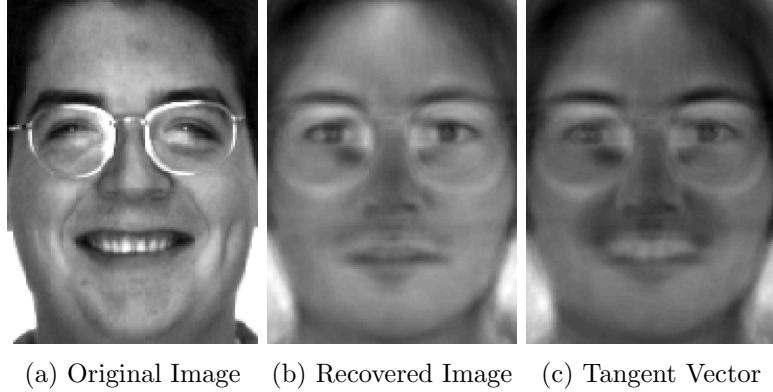


FIGURE 5.6. The tangent vector shows the subject’s smile, in particular, its shape, shading, and teeth and cheek detail (notice the right dimple), better than the recovered image. Images (b) and (c) were recovered from PCA dimension  $m_{PCA} = 30$ .

**5.4.7. Summary.** The experimental results on face databases show that LPCA-SRC can achieve higher accuracy than SRC in cases of low sampling and/or nonlinear class manifolds and small PCA feature dimension. We showed that LPCA-SRC had a significant advantage over SRC and the other algorithms for the small class sizes and nonlinear class manifolds of the AR database when the feature dimension was low. We also showed that LPCA-SRC could improve classification

## 5.4. FACE DATABASES

---

Database	$m_{PCA} = 30$		$m_{PCA} = 56$		$m_{PCA} = 120$	
	Energy	% Increased Acc. SRC / SRC <sub>pruned</sub>	Energy	% Increased Acc. SRC / SRC <sub>pruned</sub>	Energy	% Increased Acc. SRC / SRC <sub>pruned</sub>
AR-1	0.4527	3.90/3.86	0.5322	1.87/1.91	0.6522	0.80/0.60
AR-2	0.4137	3.83/2.36	0.4884	1.31/0.63	0.5988	0.62/0.53
Extended Yale B	0.3954	2.46/2.45	0.4803	1.59/1.59	0.6055	0.77/0.74
ORL	0.5385	1.34/0.05	0.6581	1.26/-0.04	0.8487	1.73/0.03

TABLE 5.8. Average energy retained in PCA dimensionality reduction (over 10 trials) to various dimensions  $m_{PCA}$  for the face database training sets, as well as the average increase in classification accuracy of LPCA-SRC over SRC and SRC<sub>pruned</sub>.

on Extended Yale B and ORL through its use of tangent vectors to provide a local approximation of the test sample and its discriminating pruning parameter, respectively.

The runtime of LPCA-SRC was sometimes much longer than that of SRC, although this was less often seen for small feature dimensions, at which LPCA-SRC tended to excel. The size of the dictionary in LPCA-SRC was observed to be a good predictor of the relationship between the runtimes of LPCA-SRC and SRC, and this could easily be computed (given estimates of the parameters  $n$  and  $d$ ) before deciding between the two methods.

To validate our claim that the tangent vectors in LPCA-SRC can contain information lost in stringent PCA dimensionality reduction, we provided examples from the AR database. We also compared the energy retained in PCA dimensionality reduction with the increase in accuracy in LPCA-SRC over SRC and saw that they were inversely correlated.

## CHAPTER 6

### Bounding the Tangent Error

In this chapter, we aim to bound the distance between a tangent vector computed in LPCA-SRC and the corresponding class manifold.

#### 6.1. Description and Assumptions

Let  $\mathcal{M}$  be a smooth  $d$ -dimensional manifold in  $\mathbb{R}^m$ , and let  $\mathbf{x}_0$  be a point on  $\mathcal{M}$ . Suppose that we use local PCA to compute the  $d$  basis vectors of the approximated tangent hyperplane of  $\mathcal{M}$  at  $\mathbf{x}_0$ . Since shifted (by  $\mathbf{x}_0$ ) and scaled versions of these basis vectors lie approximately on  $\mathcal{M}$ , this technique may be used to increase the sampling density around  $\mathbf{x}_0$ . For classification purposes, this can be used to increase the size of the training set and give us a better picture of each class manifold, as in our proposed algorithm, LPCA-SRC.

However, this approach may lead to bad training sample estimates in the sense that the tangent vectors are far from the original manifold. Local PCA may give us a poor approximation of the true tangent hyperplane under factors such as significant noise or poor sampling. It is also possible that our scaling factor produces points that are too far from  $\mathbf{x}_0$  given the local curvature.

In what follows, we assume that the scale factor  $c \geq 0$  is fixed. (Recall that  $c := r\gamma$  for  $\gamma \sim \text{unif}(0, 1)$  in LPCA-SRC.) We aim to bound, with high probability, the distance between an arbitrary LPCA-SRC tangent vector (defined as a shifted and scaled tangent hyperplane basis vector as computed using local PCA) and its closest point on  $\mathcal{M}$  for the case that  $d = 1$ . Our work uses a (modified) result of Kaslovsky and Meyer [57]. The bound we produce is stated in terms of the principal curvatures of  $\mathcal{M}$  at  $\mathbf{x}_0$ , the number and noise level of the points sampled from  $\mathcal{M}$  and used in local PCA, and the ambient dimension  $m$ . It also requires certain continuity and differentiability assumptions on  $\mathcal{M}$ .

Our work in this chapter also assumes that the local tangent hyperplane origin  $\mathbf{x}_0$  lies on the manifold  $\mathcal{M}$ . In the context of LPCA-SRC, this is not always the case, as training samples—where

## 6.2. KASLOVSKY AND MEYER'S TANGENT BOUND

---

the tangent vectors are constructed—are likely to contain noise. However, a step can be added to the offline phase of the LPCA-SRC algorithm to ensure that each local origin lies approximately on its class manifold. See Section 6.5 for details regarding this procedure.

### 6.2. Kaslovsky and Meyer's Tangent Bound

Kaslovsky and Meyer [57] computed a bound on the angle between the true tangent hyperplane and the approximated tangent hyperplane as found by local PCA. This bound holds with high probability. Let  $P := UU^\top$  be the orthogonal projector onto the true tangent space  $L$ , and let  $\hat{P} := \hat{U}\hat{U}^\top$  be the orthogonal projector onto the approximated tangent space  $\hat{L}$ . Here, the columns of  $U$  and  $\hat{U}$  are orthonormal bases of the true and approximated tangent (hyper)planes, respectively.  $U$  is unknown, but  $\hat{U}$  is found via local PCA, i.e., its columns are the first  $d$  eigenvectors of the sample covariance matrix constructed from local samples. Kaslovsky and Meyer [57] bounded (with high probability) the quantity  $\|P - \hat{P}\|_F$  in terms of (i) the manifold dimension  $d$ ; (ii) the ambient dimension  $m$ ; (iii) the number of points  $n$  in the sample covariance matrix, their noise level  $\sigma$ , and  $r$ , their maximum distance from the point  $\mathbf{x}_0$  as measured along the (true) tangent plane; and (iv) the principal curvatures of  $\mathcal{M}$  at  $\mathbf{x}_0$ . The bound also depends on two probability constants,  $\xi$  and  $\xi_\lambda$ , which can be tuned so that the probability that the bound holds is sufficiently high.

There are two conditions for the bound holding. The first, roughly, is the requirement that the true tangent space is sufficiently separable from the sample noise  $\sigma$  and the curvature of  $\mathcal{M}$  at  $\mathbf{x}_0$  at some scale (or neighborhood size) derived from  $n$  [57]. The second condition is a technical assumption, which, given that the first condition holds, will be satisfied if the sampling density is high enough [57]. Additionally, there are some “benign assumptions” on  $n$  and the probability constants  $\xi$  and  $\xi_\lambda$  [57]. We can write the Kaslovsky and Meyer bound as

$$\|P - \hat{P}\|_F \leq E(d, m, n, r, K, \sigma, \xi, \xi_\lambda),$$

where  $K$  is the  $d \times (m - d)$  matrix of principal curvatures of  $\mathcal{M}$  at  $\mathbf{x}_0$ .

Recall that in LPCA-SRC we use the version of local PCA by Singer and Wu [87], in particular, our implementation in Algorithm 4. Thus Kaslovsky and Meyer's result, derived in the case that standard local PCA is used, must be modified to hold in the Singer and Wu case. We do this in

### 6.3. MAIN THEOREM

---

Appendix A. Let us denote the modified bound (which is formally stated in Theorem A.6.2) by

$$(6.1) \quad \|P - \widehat{P}\|_F \leq E^*(d, m, n, r, K, \sigma, \xi, \xi_\lambda) =: E^*.$$

To give the reader a sense of the modified bound  $E^*$ , we can informally write

$$\|P - \widehat{P}\|_F \leq E^* \approx \frac{2\sqrt{2}}{\sqrt{n}} \frac{\left[ K^{(+)} r^3 + \sigma \sqrt{d(m-d)} \left( \sigma + \frac{r}{\sqrt{d+2}} + \frac{\mathcal{K}_W^{1/2} r^2}{2\sqrt{(d+2)(d+4)}} \right) \right]}{\frac{r^2}{d+2} - \frac{\mathcal{K}_W r^4}{4(d+2)(d+4)} - \sigma^2(\sqrt{d} + \sqrt{m-d})},$$

where  $\mathcal{K}_W$  and  $K^{(+)}$  are curvature constants derived from the entries of the curvature matrix  $K$ .

### 6.3. Main Theorem

**Notation:** Let  $\mathbf{x}_0$  be a point on the 1-dimensional manifold  $\mathcal{M}$ . Suppose that  $\mathbf{u}$  is a unit basis vector of the tangent line at  $\mathbf{x}_0$  on  $\mathcal{M}$ , and let  $\hat{\mathbf{u}}$  be the approximation of  $\mathbf{u}$  as found by Singer and Wu's local PCA [87]. For a constant  $c \geq 0$ , define  $\hat{\mathbf{v}} := c\hat{\mathbf{u}}$ , and assume without loss of generality that  $\langle \mathbf{u}, \hat{\mathbf{v}} \rangle \geq 0$ . Lastly define the (shifted and scaled) computed tangent vector  $\hat{\mathbf{z}} := \mathbf{x}_0 + \hat{\mathbf{v}}$ .

**THEOREM 6.3.1.** *Suppose that  $\mathcal{M}$  can be described by a curve  $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^m$ ,  $m \geq 2$ , with domain  $\mathcal{D} \subset \mathbb{R}$  and  $\mathbf{x}_0 = \mathbf{f}(t_0)$  for some  $t_0 \in \mathcal{D}$ , such that  $\mathbf{f}$  satisfies the following conditions:*

- (1) *The principal directions of  $\mathbf{f}$  at  $t_0$ , corresponding to the  $m-1$  principal curvatures, exist;*
- (2) *There exists  $\epsilon > 0$  such that in a neighborhood  $\mathcal{N}_\epsilon(\mathbf{x}_0) := \{\mathbf{f}(t) : \|\mathbf{f}(t) - \mathbf{f}(t_0)\|_2 < \epsilon\}$  of  $\mathbf{x}_0$ , there is a one-to-one mapping between points in  $\mathcal{N}_\epsilon(\mathbf{x}_0)$  and the tangent line of  $\mathcal{M}$  at  $\mathbf{x}_0$ . Further,  $\mathbf{f}$  is parametrized so that  $I_\epsilon(\mathbf{x}_0) := \{t \in \mathcal{D} : \mathbf{f}(t) \in \mathcal{N}_\epsilon(\mathbf{x}_0)\}$  is an interval on which  $\mathbf{f}$  is continuous and one-to-one;*
- (3) *The scaling factor  $c$  is small enough so that there exists  $\mathbf{f}(t_c) \in \mathcal{N}_\epsilon(\mathbf{x}_0)$  that is mapped to the point  $\mathbf{x}_0 + c\mathbf{u}$  on the tangent line via the one-to-one mapping in (2);*
- (4)  *$\mathbf{f}$  is three times differentiable on  $I_\epsilon(\mathbf{x}_0)$  with continuous second-order derivatives in  $\overline{I_\epsilon(\mathbf{x}_0)}$ , and the third derivative of  $\mathbf{f}$  satisfies  $\sup_{\mathbf{f}(t) \in \mathcal{N}_\epsilon(\mathbf{x}_0)} \|\mathbf{f}'''(t)\|_2 =: M < \infty$ .*

Suppose that all the conditions of the modified Kaslovsky and Meyer tangent bound [57] hold so that  $\|P - \widehat{P}\|_F \leq E^*$  with probability greater than  $p := 1 - 2e^{-\xi_\lambda^2} - 9e^{-\xi^2}$ . Then the distance

## 6.4. PROOF OF MAIN THEOREM

---

between the shifted and scaled computed tangent vector  $\hat{\mathbf{z}}$  and the manifold  $\mathcal{M}$  satisfies

$$(6.2) \quad \text{dist}(\hat{\mathbf{z}} - \mathcal{M}) := \inf_{\mathbf{z}' \in \mathcal{M}} \|\hat{\mathbf{z}} - \mathbf{z}'\|_2 \leq c \sqrt{\frac{E^*}{2}} + \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R} \right)^{1/2} =: \delta$$

with probability greater than  $p$ , where

$$|\mathcal{R}| \leq R_{\text{bound}} := \frac{Mc^5}{6}(m-1) \left( \kappa + \frac{Mc}{6} \right).$$

Here,  $\kappa$  is given by

$$\kappa := \max_{i \in \{2, \dots, m\}} |\kappa^{(i)}|,$$

where  $\kappa^{(i)}$  is the principal curvature in the  $i$ th normal direction of  $\mathcal{M}$  at  $\mathbf{x}_0$ ,  $i \in \{2, \dots, m\}$ .

Set  $\boldsymbol{\kappa} = [\kappa^{(2)}, \dots, \kappa^{(m)}]^T \in \mathbb{R}^{m-1}$  and define  $\delta_{\text{approx}} := c \sqrt{\frac{E^*}{2}} + \frac{c^2}{2} \|\boldsymbol{\kappa}\|_2$  and

$$\zeta_{\text{bound}} := \max \left\{ \left( \frac{c^4}{4} \|\boldsymbol{\kappa}\|_2^2 + R_{\text{bound}} \right)^{1/2} - \frac{c^2}{2} \|\boldsymbol{\kappa}\|_2, \frac{c^2}{2} \|\boldsymbol{\kappa}\|_2 - \left( \frac{c^4}{4} \|\boldsymbol{\kappa}\|_2^2 - R_{\text{bound}} \right)^{1/2} \right\}$$

Then  $\delta_{\text{approx}}$  approximates the upper bound  $\delta \geq \text{dist}(\hat{\mathbf{z}} - \mathcal{M})$  with  $|\delta - \delta_{\text{approx}}| \leq \zeta_{\text{bound}}$  with probability greater than  $p$ .

### 6.4. Proof of Main Theorem

We prove the main part of Theorem 6.3.1 using two lemmas, each of which essentially restates half the theorem.

#### 6.4.1. First Lemma.

LEMMA 6.4.1. Let  $\mathcal{M}$  be a 1-dimensional manifold in  $\mathbb{R}^m$  ( $m \geq 2$ ) with  $\mathbf{x}_0 \in \mathcal{M}$ , and suppose that  $\mathbf{u}$  is a unit basis vector of the tangent line at  $\mathbf{x}_0$  on  $\mathcal{M}$ . For fixed  $c \geq 0$ , obtain a scaled and shifted version of  $\mathbf{u}$  by setting  $\mathbf{z} := \mathbf{x}_0 + c\mathbf{u}$ . Suppose that  $\mathcal{M}$  can be described by a curve  $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^m$ , with domain  $\mathcal{D} \subset \mathbb{R}$  and  $\mathbf{x}_0 = \mathbf{f}(t_0)$  for some  $t_0 \in \mathcal{D}$ , such that the properties (1)-(4) in Theorem 6.3.1 hold. Then the distance between  $\mathbf{z}$  and  $\mathcal{M}$  satisfies

$$\text{dist}(\mathbf{z} - \mathcal{M}) \leq \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R},$$

## 6.4. PROOF OF MAIN THEOREM

---

where

$$|\mathcal{R}| \leq \frac{Mc^5}{6}(m-1) \left( \kappa + \frac{Mc}{6} \right).$$

Here,  $\kappa$  is given by

$$\kappa := \max_{i \in \{2, \dots, m\}} |\kappa^{(i)}|,$$

where  $\kappa^{(i)}$  is the principal curvature in the  $i$ th normal direction of  $\mathcal{M}$  at  $\mathbf{x}_0$ ,  $i \in \{2, \dots, m\}$ .

**Proof of Lemma 6.4.1.** The proof of this lemma relies heavily on the geometric data model in the work of Kaslovsky and Meyer [57] (see also the work of Tyagi et al. [94]).

*Geometric setup in three dimensions.* Let  $m = 3$ . Suppose that  $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^3$  is a smooth curve with coordinate functions  $f_i : \mathbb{R} \rightarrow \mathbb{R}$ ,  $i = 1, 2, 3$ , and point on the curve

$$\mathbf{f}(t_0) = [f_1(t_0), f_2(t_0), f_3(t_0)]^\top =: \mathbf{x}_0.$$

By the assumptions in Lemma 6.4.1, we can compute the unit tangent vector  $\mathbf{t}_{t_0}$ , the unit normal vector  $\mathbf{n}_{t_0}$ , and the unit binormal vector  $\mathbf{b}_{t_0}$  at  $t_0$  on  $\mathbf{f}$ . By definition, all three vectors are mutually orthogonal. Further,  $\mathbf{n}_{t_0}$  indicates the normal direction of curvature, and similarly,  $\mathbf{b}_{t_0}$  indicates the normal direction of torsion, or second curvature.

Let  $R$  be the rotation matrix that maps  $\mathbf{n}_{t_0}$  to  $\mathbf{e}_2 := [0, 1, 0]^\top$  and  $\mathbf{b}_{t_0}$  to  $\mathbf{e}_3 := [0, 0, 1]^\top$ . We define the function  $\tilde{\mathbf{f}} := R(\mathbf{f} - \mathbf{x}_0)$ . This transformation shifts  $\mathbf{x}_0$  to the origin and aligns the normal directions of curvature (i.e., the principal directions) of  $\mathbf{f}$  at  $t_0$  with the last two coordinate axes. As a result of the orthogonality relationships, the unit tangent vector becomes aligned with the  $x$ -axis, and since  $\mathbf{b}_{t_0} := \mathbf{t}_{t_0} \times \mathbf{n}_{t_0}$ ,  $R(\mathbf{t}_{t_0}) = \mathbf{e}_1 := [1, 0, 0]^\top$ . The function  $\tilde{\mathbf{f}}$  can subsequently be parametrized locally by  $\tilde{\mathbf{f}}(x) = [x, \tilde{f}_2(x), \tilde{f}_3(x)]^\top$ , where  $\tilde{f}_i : \mathbb{R} \rightarrow \mathbb{R}$ ,  $i = 2, 3$  [69, 94], as follows:

By the assumptions in Lemma 6.4.1, there exists  $\epsilon$  such that each point  $\mathbf{f}(t) \in \mathcal{N}_\epsilon(\mathbf{x}_0)$  has a unique mapping onto the tangent line at  $\mathbf{x}_0$ , and so there exists a bijective function  $\gamma : I_\epsilon(\mathbf{x}_0) \rightarrow [-\tilde{\epsilon}, \tilde{\epsilon}]$  for some  $0 < \tilde{\epsilon} \leq \epsilon$  (since we assumed that there is exactly one point  $\mathbf{f}(t) \in \mathcal{N}_\epsilon(\mathbf{x}_0)$  for each  $t \in I_\epsilon(\mathbf{x}_0)$ ). Further, since  $\mathbf{f}$  is continuous on the interval  $I_\epsilon(\mathbf{x}_0)$ , we can define  $x := \gamma(t) = \pm(t - t_0)$ ,

## 6.4. PROOF OF MAIN THEOREM

---

which renders the desired parametrization of  $\tilde{\mathbf{f}}$  as

$$\tilde{\mathbf{f}}(x) = \tilde{\mathbf{f}}(\gamma(t)) := R(\mathbf{f}(t) - \mathbf{x}_0),$$

for  $\mathbf{f}(t)$  close to  $\mathbf{f}(t_0)$ . It follows that the point  $\mathbf{f}(t_0) = \mathbf{x}_0$  corresponds to  $\tilde{\mathbf{f}}(0) = \mathbf{0}$  via

$$\mathbf{0} = R(0) = R(\mathbf{f}(t_0) - \mathbf{x}_0) = \tilde{\mathbf{f}}(\gamma(t_0)) = \tilde{\mathbf{f}}(\pm(t_0 - t_0)) = \tilde{\mathbf{f}}(0).$$

To further understand this correspondence, let  $\mathbf{z} \in \mathbb{R}^3$  be a point on the tangent line of  $\mathbf{f}$  at  $t_0$  such that  $\|\mathbf{z} - \mathbf{x}_0\|_2 = c$ , for small  $c \geq 0$  (i.e.,  $c \leq \tilde{\epsilon}$ ). There are two such points, and they can be differentiated by setting  $\mathbf{z}_+$  to be the point in the direction of the unit tangent vector  $\mathbf{t}_{t_0}$  and  $\mathbf{z}_-$  to be the point in the opposite direction. Then  $\mathbf{z}_+$  corresponds to  $[c, 0, 0]^\top$  and  $\mathbf{z}_-$  corresponds to  $[-c, 0, 0]^\top$  in the shifted and rotated coordinate system. Further, there is a unique point  $\mathbf{x}_{\mathbf{z}_+} = \mathbf{f}(t_{\mathbf{z}_+})$  in the original coordinate system that corresponds to  $\tilde{\mathbf{f}}(c)$ , where  $c = \gamma(t_{\mathbf{z}_+}) = \pm(t_{\mathbf{z}_+} - t_0)$  and  $\mathbf{f}(t_{\mathbf{z}_+})$  is close to  $\mathbf{f}(t_0)$  (i.e.,  $\|\mathbf{f}(t_{\mathbf{z}_+}) - \mathbf{f}(t_0)\|_2 \leq \epsilon$ ). The precise correspondence is

$$\tilde{\mathbf{f}}(c) = \tilde{\mathbf{f}}(\gamma(t_{\mathbf{z}_+})) = R(\mathbf{f}(t_{\mathbf{z}_+}) - \mathbf{x}_0) = R(\mathbf{x}_{\mathbf{z}_+} - \mathbf{x}_0).$$

An analogous statement holds for  $\mathbf{z}_-$ .

The advantages of viewing points in the shifted and rotated coordinate system (e.g.,  $[c, 0, 0]^\top$  instead of  $\mathbf{z}_+$ ) are (i) the relationship between the function and the tangent line are preserved, so that we may use  $\tilde{\mathbf{f}}$  instead of  $\mathbf{f}$  without loss of generality to measure the distance between a point near  $\mathbf{x}_0$  on the tangent line and the curve  $\mathbf{f}$ , and (ii) the first nonzero terms in the Taylor expansion of  $\tilde{f}_2$  and  $\tilde{f}_3$  at 0 are quadratic:

$$\begin{aligned}\tilde{f}_i(x) &= \frac{1}{2}\tilde{f}_i''(x)x^2 + \text{higher-order terms} \\ &= \frac{1}{2}\kappa^{(i)}x^2 + \text{higher-order terms},\end{aligned}$$

for  $i = 2, 3$  [45]. Here,  $\tilde{f}_i''$  denotes the second derivative of  $\tilde{f}_i$ ,  $i = 2, 3$ . Additionally,  $\kappa^{(2)}$  is the curvature of  $\tilde{\mathbf{f}}$  at 0, and  $\kappa^{(3)}$  is the torsion of  $\tilde{\mathbf{f}}$  at 0; equivalently,  $\kappa^{(i)}$  is the curvature in the  $i$ th

## 6.4. PROOF OF MAIN THEOREM

---

normal direction,  $i = 2, 3$ . Since curvature is invariant under rigid transformation, these values can be obtained from the original function  $\mathbf{f}$  at  $t_0$ . More likely, however, since an explicit expression for  $\mathbf{f}$  may be difficult or impossible to obtain in practice, these values can be estimated directly from  $\mathcal{M}$ .

*Geometric setup in arbitrary ambient dimension.* This geometric setup is easily extended for general  $m$ . The curvature vectors  $\mathbf{n}_{t_0}$  and  $\mathbf{b}_{t_0}$  are replaced with the principal directions, i.e., (unit) vectors indicating the  $m - 1$  normal directions of principal curvature. We can write these as  $\mathbf{n}_{t_0,2}, \dots, \mathbf{n}_{t_0,m}$ . Together with  $\mathbf{t}_{t_0}$ , these vectors are mutually orthogonal. We set  $R$  to be the rotation matrix that maps  $\mathbf{n}_{t_0,2}, \dots, \mathbf{n}_{t_0,m}$  to  $\mathbf{e}_2, \dots, \mathbf{e}_m$ , respectively, and we once again set  $\tilde{\mathbf{f}} = R(\mathbf{f} - \mathbf{x}_0)$ . Note that the principal directions and  $R$  should be chosen so that  $R(\mathbf{t}_{t_0}) = \mathbf{e}_1$  (instead of  $-\mathbf{e}_1$ ), for consistency.

The function  $\tilde{\mathbf{f}} : \mathbb{R} \rightarrow \mathbb{R}^m$  can be described locally by  $\tilde{\mathbf{f}}(x) = [x, \tilde{f}_2(x), \dots, \tilde{f}_m(x)]^\top$ , with

$$\begin{aligned}\tilde{f}_i(x) &= \frac{1}{2} \tilde{f}_i''(x) x^2 + \text{higher-order terms} \\ &= \frac{1}{2} \kappa^{(i)} x^2 + \text{higher-order terms},\end{aligned}$$

where  $\kappa^{(i)}$  is the curvature of  $\tilde{\mathbf{f}}$  at 0 (or equivalently,  $\mathbf{f}$  at  $t_0$ ) in the  $i$ th normal direction,  $i = 2, \dots, m$ .

*Leveraging the geometric setup to prove Lemma 6.4.1.* Recall that our goal is to bound the distance between  $\mathbf{z} := \mathbf{x}_0 + c\mathbf{u}$  and  $\mathcal{M}$ . We assume that  $\|\mathbf{u}\|_2 = 1$  so that  $\mathbf{u} = \pm\mathbf{t}_{t_0}$ . The point corresponding to  $\mathbf{z}$  in the rotated and shifted coordinate system is  $[x_{\mathbf{z}}, 0, \dots, 0]$ , where  $x_{\mathbf{z}} = c \langle \mathbf{u}, \mathbf{t}_{t_0} \rangle = c \operatorname{sgn}(\langle \mathbf{u}, \mathbf{t}_{t_0} \rangle)$ .

By the geometric setup, we have that

$$\operatorname{dist}(\mathbf{z} - \mathcal{M}) = \min_{t \in \mathcal{D}} \|\mathbf{z} - \mathbf{f}(t)\|_2 \leq \min_{x \in [-\tilde{\epsilon}, \tilde{\epsilon}]} \| [x_{\mathbf{z}}, 0, \dots, 0]^\top - \tilde{\mathbf{f}}(x) \|_2.$$

This, along with the observation that  $\min_{x \in [-\tilde{\epsilon}, \tilde{\epsilon}]} \| [x_{\mathbf{z}}, 0, \dots, 0]^\top - \tilde{\mathbf{f}}(x) \|_2 \leq \| [x_{\mathbf{z}}, 0, \dots, 0]^\top - \tilde{\mathbf{f}}(x_{\mathbf{z}}) \|_2$ , yields

$$\begin{aligned}\operatorname{dist}(\mathbf{z} - \mathcal{M}) &\leq \| [x_{\mathbf{z}}, 0, \dots, 0]^\top - \tilde{\mathbf{f}}(x_{\mathbf{z}}) \|_2 \\ &= \| [x_{\mathbf{z}}, 0, \dots, 0]^\top - [x_{\mathbf{z}}, \tilde{f}_2(x_{\mathbf{z}}), \dots, \tilde{f}_m(x_{\mathbf{z}})]^\top \|_2\end{aligned}$$

## 6.4. PROOF OF MAIN THEOREM

---

$$\begin{aligned}
&= \left\| [\tilde{f}_2(x_{\mathbf{z}}), \dots, \tilde{f}_m(x_{\mathbf{z}})]^T \right\|_2 \\
&\approx \left\| \left[ \frac{1}{2} \kappa^{(2)} x_{\mathbf{z}}^2, \dots, \frac{1}{2} \kappa^{(m)} x_{\mathbf{z}}^2 \right]^T \right\|_2 \\
&= \frac{x_{\mathbf{z}}^2}{2} \left( \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} \\
&= \frac{c^2}{2} \left( \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2},
\end{aligned}$$

where the approximation follows from Taylor's expansion of the functions  $\tilde{f}_i$ ,  $i = 2, \dots, m$ , around the origin. It is clear from this sequence of equations that the distance between a point on the tangent line at  $\mathbf{x}_0$  and the manifold  $\mathcal{M}$  depends primarily on the scaling factor  $c$  (how close it is to  $\mathbf{x}_0$  on the tangent line) and the principal curvatures in the  $m - 1$  normal directions.

Define the interval

$$I_{x_{\mathbf{z}}} := \begin{cases} (0, x_{\mathbf{z}}), & x_{\mathbf{z}} \geq 0, \\ (x_{\mathbf{z}}, 0), & x_{\mathbf{z}} < 0. \end{cases}$$

If we further assume that the third derivative  $\tilde{f}_i'''(x)$  exists for all  $x \in I_{x_{\mathbf{z}}}$  and that the second derivative  $\tilde{f}_i''(x)$  is continuous for all  $x \in \bar{I}_{x_{\mathbf{z}}}$ ,  $i = 2, \dots, m$ , then Taylor's theorem states that

$$\begin{aligned}
\|\mathbf{z} - \mathcal{M}\|_2 &\leq \left\| [\tilde{f}_2(x_{\mathbf{z}}), \dots, \tilde{f}_m(x_{\mathbf{z}})]^T \right\|_2 \\
&= \left\| \left[ \frac{1}{2} \kappa^{(2)} x_{\mathbf{z}}^2 + \frac{\tilde{f}_2'''(\xi_2)}{3!} x_{\mathbf{z}}^3, \dots, \frac{1}{2} \kappa^{(m)} x_{\mathbf{z}}^2 + \frac{\tilde{f}_m'''(\xi_m)}{3!} x_{\mathbf{z}}^3 \right]^T \right\|_2 \\
&= \left( \sum_{i=2}^m \left( \frac{1}{2} \kappa^{(i)} x_{\mathbf{z}}^2 + \frac{\tilde{f}_i'''(\xi_i)}{3!} x_{\mathbf{z}}^3 \right)^2 \right)^{1/2} \\
&= \left( \sum_{i=2}^m \left( \frac{1}{4} (\kappa^{(i)})^2 x_{\mathbf{z}}^4 + \frac{1}{6} \kappa^{(i)} \tilde{f}_i'''(\xi_i) x_{\mathbf{z}}^5 + \frac{1}{36} (\tilde{f}_i'''(\xi_i))^2 x_{\mathbf{z}}^6 \right) \right)^{1/2} \\
&= \left( \frac{x_{\mathbf{z}}^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \frac{x_{\mathbf{z}}^5}{6} \sum_{i=2}^m \tilde{f}_i'''(\xi_i) \left( \kappa^{(i)} + \frac{x_{\mathbf{z}}}{6} \tilde{f}_i'''(\xi_i) \right) \right)^{1/2} \\
&= \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \frac{x_{\mathbf{z}}^5}{6} \sum_{i=2}^m \tilde{f}_i'''(\xi_i) \left( \kappa^{(i)} + \frac{x_{\mathbf{z}}}{6} \tilde{f}_i'''(\xi_i) \right) \right)^{1/2},
\end{aligned}$$

## 6.4. PROOF OF MAIN THEOREM

---

where each  $\xi_i$  is some number in  $I_{x_z}$ ,  $i = 2, \dots, m$ . For simplicity, we write

$$\|[\tilde{f}_2(x_z), \dots, \tilde{f}_m(x_z)]^\top\|_2^2 = \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R},$$

where

$$\mathcal{R} := \frac{x_z^5}{6} \sum_{i=2}^m \tilde{f}_i'''(\xi_i) \left( \kappa^{(i)} + \frac{x_z}{6} \tilde{f}_i'''(\xi_i) \right).$$

We want to bound  $|\mathcal{R}|$ . Since  $I_{x_z} \subset [-\tilde{\epsilon}, \tilde{\epsilon}]$  when  $c \leq \tilde{\epsilon}$ , then for any  $x \in I_{x_z}$ , there exists  $t \in I_\epsilon(\mathbf{x}_0)$  such that  $x = \gamma(t) = \pm(t - t_0)$ . Because  $\mathbf{f}(t) = [f_1(t), \dots, f_m(t)]^\top$  and  $\tilde{\mathbf{f}}(x) = R(\mathbf{f}(t) - \mathbf{x}_0)$ , then if  $R_{ij}$  denotes the  $(i, j)$ th entry of the rotation matrix  $R$ , it follows that

$$\begin{aligned} |\tilde{f}_i'''(x)| &= |R_{i1}f_1'''(t) + \dots + R_{im}f_m'''(t)| \\ &= |\langle [R_{i1}, \dots, R_{im}]^\top, \mathbf{f}'''(t) \rangle| \\ &\leq \|\mathbf{f}'''(t)\|_2 \end{aligned}$$

by Cauchy-Schwarz, since  $R$  is an orthogonal matrix,  $\gamma'(t) = \pm 1$ , and  $\gamma''(t) = \gamma'''(t) = 0$ . So if we set

$$M := \max_{t \in I_\epsilon(\mathbf{x}_0)} \|\mathbf{f}'''(t)\|_2 = \max_{\mathbf{f}(t) \in \mathcal{N}_\epsilon(\mathbf{x}_0)} \|\mathbf{f}'''(t)\|_2 \quad \text{and} \quad \kappa := \max_{i \in \{2, \dots, m\}} |\kappa^{(i)}|,$$

it follows that

$$\begin{aligned} |\mathcal{R}| &= \left| \frac{x_z^5}{6} \sum_{i=2}^m \tilde{f}_i'''(\xi_i) \left( \kappa^{(i)} + \frac{x_z}{6} \tilde{f}_i'''(\xi_i) \right) \right| \\ &= \left| \frac{x_z^5}{6} \left| \sum_{i=2}^m \tilde{f}_i'''(\xi_i) \left( \kappa^{(i)} + \frac{x_z}{6} \tilde{f}_i'''(\xi_i) \right) \right| \right| \\ &\leq \left| \frac{(c \langle \mathbf{u}, \mathbf{t}_{t_0} \rangle)^5}{6} \right| \left[ \sum_{i=2}^m \left| \tilde{f}_i'''(\xi_i) \right| \left( |\kappa^{(i)}| + \left| \frac{c \langle \mathbf{u}, \mathbf{t}_{t_0} \rangle}{6} \right| \left| \tilde{f}_i'''(\xi_i) \right| \right) \right] \\ &\leq \frac{c^5}{6} \left[ \sum_{i=2}^m \left| \tilde{f}_i'''(\xi_i) \right| \left( |\kappa^{(i)}| + \frac{c}{6} \left| \tilde{f}_i'''(\xi_i) \right| \right) \right] \\ &\leq \frac{c^5}{6} \left[ \sum_{i=2}^m M \left( |\kappa^{(i)}| + \frac{c}{6} M \right) \right] \end{aligned}$$

## 6.4. PROOF OF MAIN THEOREM

---

$$\leq \frac{Mc^5}{6}(m-1)\left(\kappa + \frac{Mc}{6}\right).$$

This is sufficient to prove Lemma 6.4.1. The differentiability and continuity assumptions on  $\mathbf{f}$  ensure that  $\tilde{\mathbf{f}}$  is three times differentiable on  $(-c, c)$  with continuous second-order derivatives on  $[-c, c]$ , so we can apply Taylor's theorem.  $\square$

**6.4.2. Second Lemma.** Again, recall that in LPCA-SRC we use the version of local PCA by Singer and Wu [87]. Thus Kaslovsky and Meyer's result [57] must be modified to hold in this case. We do this in Appendix A, and we denote the modified bound using Eq. (6.1).

LEMMA 6.4.2. *Let  $\mathbf{x}_0$  be a point on the 1-dimensional manifold  $\mathcal{M}$ . Define  $\mathbf{u}$  to be a unit basis vector of the tangent line at  $\mathbf{x}_0$  on  $\mathcal{M}$ , and let  $\hat{\mathbf{u}}$  be the approximation of  $\mathbf{u}$  as found by Singer and Wu's local PCA [87]. Suppose that all the conditions of Kaslovsky and Meyer's modified tangent bound [57] hold so that  $\|P - \hat{P}\|_F \leq E^*$  with probability greater than  $p := 1 - 2e^{-\xi_\lambda^2} - 9e^{-\xi^2}$ . Then*

$$\|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2 \leq \sqrt{\frac{E^*}{2}}$$

*with probability greater than  $p$ , where  $P_{\mathbf{u}}(\hat{\mathbf{u}})$  is the orthogonal projection of  $\hat{\mathbf{u}}$  onto the subspace spanned by  $\mathbf{u}$  (i.e., onto the tangent line at  $\mathbf{x}_0$ ).*

**Proof of Lemma 6.4.2.** Since  $d = 1$ ,  $\|P - \hat{P}\|_F = \|\mathbf{u}\mathbf{u}^\top - \hat{\mathbf{u}}\hat{\mathbf{u}}^\top\|_F$ . We can assume without loss of generality that  $\langle \mathbf{u}, \hat{\mathbf{u}} \rangle \geq 0$ , since  $\mathbf{u}\mathbf{u}^\top = (-\mathbf{u})(-\mathbf{u})^\top$ . Let the  $(i, j)$ th entry of  $P := \mathbf{u}\mathbf{u}^\top \in \mathbb{R}^{m \times m}$  be given by  $P_{ij} = u_i u_j$ , where  $u_i$  denotes the  $i$ th coordinate of  $\mathbf{u}$ ,  $1 \leq i \leq m$ . It follows that the Frobenius norm of the difference between  $P$  and  $\hat{P}$  satisfies

$$\begin{aligned} \|P - \hat{P}\|_F^2 &= \sum_{i=1}^m \sum_{j=1}^m (P_{ij} - \hat{P}_{ij})^2 \\ &= \sum_{i=1}^m \sum_{j=1}^m (u_i u_j - \hat{u}_i \hat{u}_j)^2 \\ &= \sum_{i=1}^m u_i^2 \sum_{j=1}^m u_j^2 + \sum_{i=1}^m \hat{u}_i^2 \sum_{j=1}^m \hat{u}_j^2 - 2 \sum_{i=1}^m u_i \hat{u}_i \sum_{j=1}^m u_j \hat{u}_j \\ &= 2 - 2 \langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2. \end{aligned}$$

## 6.4. PROOF OF MAIN THEOREM

---

The modified Kaslovsky and Meyer bound [57] states that  $\|P - \widehat{P}\|_F^2 = 2 - 2\langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \leq E^*$ .

*Case 1:*  $E^* > 2$ . Of course,  $2 - 2\langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \leq 2$ . So if  $E^* > 2$ , the modified Kaslovsky and Meyer bound [57] is superfluous. In this case, it suffices to show that  $\|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2^2 \leq 1$ , since  $1 < \frac{E^*}{2}$ . We have

$$\begin{aligned} \|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2^2 &= \|\hat{\mathbf{u}} - \langle \mathbf{u}, \hat{\mathbf{u}} \rangle \mathbf{u}\|_2^2 \\ &= \langle \hat{\mathbf{u}}, \hat{\mathbf{u}} \rangle - 2\langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 + \langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \langle \mathbf{u}, \mathbf{u} \rangle \\ &= 1 - \langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \leq 1. \quad \square \end{aligned}$$

*Case 2:*  $E^* \leq 2$ . In this case,  $1 - \frac{E^*}{2} \geq 0$ , and we have

$$\begin{aligned} \|P - \widehat{P}\|_F^2 = 2 - 2\langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \leq E^* &\Rightarrow 1 - \frac{E^*}{2} \leq \langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \\ &\Rightarrow \sqrt{1 - \frac{E^*}{2}} \leq \langle \mathbf{u}, \hat{\mathbf{u}} \rangle \\ &\Rightarrow \sqrt{1 - \frac{E^*}{2}} \leq \cos(\theta), \end{aligned}$$

where  $\theta$  is the angle between  $\mathbf{u}$  and  $\hat{\mathbf{u}}$ . The second-to-last inequality follows from  $\langle \mathbf{u}, \hat{\mathbf{u}} \rangle \geq 0$ .

Consider the right triangle with angle  $\theta$  placed at the origin and incident sides  $\hat{\mathbf{u}}$  and  $P_{\mathbf{u}}(\hat{\mathbf{u}})$ . We want to bound the length of the third side, which is given by  $\|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2$ . Since  $\|\hat{\mathbf{u}}\|_2 = 1$ , the above inequality and basic trigonometry implies

$$(6.3) \quad \sqrt{1 - \frac{E^*}{2}} \leq \cos(\theta) = \|P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2.$$

By the Pythagorean theorem,

$$(6.4) \quad \|P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2^2 + \|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2^2 = 1.$$

Combining Eq. (6.3) and Eq. (6.4) produces

$$\begin{aligned} 1 - \frac{E^*}{2} + \|P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2^2 &\leq 1 \\ \Rightarrow \|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2^2 &\leq 1 - \left(1 - \frac{E^*}{2}\right) \end{aligned}$$

## 6.4. PROOF OF MAIN THEOREM

---

$$\Rightarrow \|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{u}})\|_2^2 \leq \frac{E^*}{2}. \quad \square$$

COROLLARY 6.4.1. *Under the same hypotheses as Lemma 6.4.2, if we replace  $\hat{\mathbf{u}}$  with its scaled version  $\hat{\mathbf{v}} := c\hat{\mathbf{u}}$  for  $c \geq 0$ , then*

$$\|\hat{\mathbf{v}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2^2 \leq \frac{c^2 E^*}{2}$$

*with probability greater than  $p$ .*

**Proof of Corollary 6.4.1.** We again assume without loss of generality that  $\langle \mathbf{u}, \hat{\mathbf{u}} \rangle > 0$  so that  $\langle \mathbf{u}, \hat{\mathbf{v}} \rangle \geq 0$ .

*Case 1:  $E^* > 2$ .* As we saw above, the modified Kaslovsky and Meyer bound [57] is unhelpful, since

$$\|P - \hat{P}\|_F^2 = 2 - 2 \langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \leq 2.$$

The corollary still holds, however, as

$$\begin{aligned} \|\hat{\mathbf{v}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2^2 &= \|\hat{\mathbf{v}} - \langle \mathbf{u}, \hat{\mathbf{v}} \rangle \mathbf{u}\|_2^2 = \langle \hat{\mathbf{v}}, \hat{\mathbf{v}} \rangle - 2 \langle \mathbf{u}, \hat{\mathbf{v}} \rangle^2 + \langle \mathbf{u}, \hat{\mathbf{v}} \rangle^2 \langle \mathbf{u}, \mathbf{u} \rangle \\ &= \langle c\hat{\mathbf{u}}, c\hat{\mathbf{u}} \rangle - \langle \mathbf{u}, c\hat{\mathbf{u}} \rangle^2 \\ &= c^2 - c^2 \langle \mathbf{u}, \hat{\mathbf{u}} \rangle^2 \\ &\leq c^2 < \frac{c^2 E^*}{2}. \quad \square \end{aligned}$$

*Case 2:  $E^* \leq 2$ .* In this case,  $1 - \frac{E^*}{2} \geq 0$ , and as we saw before,  $\sqrt{1 - \frac{E^*}{2}} \leq \cos(\theta)$ . We previously defined  $\theta$  to be the angle between  $\mathbf{u}$  and  $\hat{\mathbf{u}}$ , but of course, it is also the angle between  $\mathbf{u}$  and  $\hat{\mathbf{v}} = c\hat{\mathbf{u}}$ . Similarly to above, we consider the right triangle formed by the vectors  $\hat{\mathbf{v}}$  and  $P_{\mathbf{u}}(\hat{\mathbf{v}})$ . We want to bound the length of the third side,  $\|\hat{\mathbf{v}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2$ . Since  $\|\hat{\mathbf{v}}\|_2 = c$ , it follows that

$$\begin{aligned} \sqrt{1 - \frac{E^*}{2}} &\leq \cos(\theta) = \frac{\|P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2}{c} \\ \Rightarrow c^2 \left(1 - \frac{E^*}{2}\right) &\leq \|P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2^2. \end{aligned}$$

## 6.4. PROOF OF MAIN THEOREM

---

By the Pythagorean theorem,

$$\|P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2^2 + \|\hat{\mathbf{v}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2^2 = c^2,$$

and combining these equations produces

$$\begin{aligned} c^2 \left(1 - \frac{E^*}{2}\right) + \|\hat{\mathbf{u}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2^2 &\leq c^2 \\ \Rightarrow \|\hat{\mathbf{v}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2^2 &\leq c^2 - c^2 \left(1 - \frac{E^*}{2}\right) \leq \frac{c^2 E^*}{2}. \quad \square \end{aligned}$$

**6.4.3. Combining the Lemmas to Prove the Main Part of Theorem 6.3.1.** We first review the notation and make a new definition. We have that  $\mathbf{u}$  is a unit basis vector of the tangent line at  $\mathbf{x}_0$  on  $\mathcal{M}$  and  $\hat{\mathbf{u}}$  is the approximation of  $\mathbf{u}$  as found by Singer and Wu's local PCA [87]. Without loss of generality, we assume that  $\langle \mathbf{u}, \hat{\mathbf{u}} \rangle \geq 0$ . For scalar  $c \geq 0$ , we set  $\hat{\mathbf{v}} := c\hat{\mathbf{u}}$ . We can write

$$P_{\mathbf{u}}(\hat{\mathbf{v}}) = c \langle \mathbf{u}, \hat{\mathbf{u}} \rangle \mathbf{u} = \tilde{c} \mathbf{u},$$

where  $\tilde{c} := c \langle \mathbf{u}, \hat{\mathbf{u}} \rangle \geq 0$ . Set  $\mathbf{v} := \tilde{c} \mathbf{u}$ . Further, (re)define  $\mathbf{z} := \mathbf{x}_0 + \mathbf{v}$  and  $\hat{\mathbf{z}} := \mathbf{x}_0 + \hat{\mathbf{v}}$ .

To prove the main part of Theorem 6.3.1, let  $\mathcal{A}$  be the affine subspace given by  $\mathcal{A} := \{\mathbf{x}_0 + \alpha \mathbf{u} : \alpha \in \mathbb{R}\}$ . Then

$$(6.5) \quad \text{dist}(\hat{\mathbf{z}} - \mathcal{M}) \leq \|\hat{\mathbf{z}} - P_{\mathcal{A}}(\hat{\mathbf{z}})\|_2 + \text{dist}(P_{\mathcal{A}}(\hat{\mathbf{z}}) - \mathcal{M}).$$

Here,  $P_{\mathcal{A}}(\hat{\mathbf{z}})$  is the orthogonal projection of  $\hat{\mathbf{z}}$  onto  $\mathcal{A}$ , which can be written as  $P_{\mathcal{A}}(\hat{\mathbf{z}}) = \mathbf{x}_0 + P_{\mathbf{u}}(\hat{\mathbf{z}} - \mathbf{x}_0)$ . Observe that

$$\begin{aligned} \|\hat{\mathbf{z}} - P_{\mathcal{A}}(\hat{\mathbf{z}})\|_2 &= \|\hat{\mathbf{z}} - (\mathbf{x}_0 + P_{\mathbf{u}}(\hat{\mathbf{z}} - \mathbf{x}_0))\|_2 \\ &= \|(\mathbf{x}_0 + \hat{\mathbf{v}}) - \mathbf{x}_0 - P_{\mathbf{u}}(\hat{\mathbf{z}} - \mathbf{x}_0)\|_2 \\ &= \|\hat{\mathbf{v}} - P_{\mathbf{u}}((\mathbf{x}_0 + \hat{\mathbf{v}}) - \mathbf{x}_0)\|_2 \\ &= \|\hat{\mathbf{v}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2, \end{aligned}$$

## 6.4. PROOF OF MAIN THEOREM

---

and

$$\begin{aligned}
\text{dist}(P_{\mathcal{A}}(\hat{\mathbf{z}}) - \mathcal{M}) &= \text{dist}((\mathbf{x}_0 + P_{\mathbf{u}}(\hat{\mathbf{v}})) - \mathcal{M}) \\
&= \text{dist}((\mathbf{x}_0 + \langle \mathbf{u}, \hat{\mathbf{v}} \rangle \mathbf{u}) - \mathcal{M}) \\
&= \text{dist}((\mathbf{x}_0 + c_1 \langle \mathbf{u}, \hat{\mathbf{u}} \rangle \mathbf{u}) - \mathcal{M}) \\
&= \text{dist}((\mathbf{x}_0 + \mathbf{v}) - \mathcal{M}) \\
&= \text{dist}(\mathbf{z} - \mathcal{M}).
\end{aligned}$$

Plugging these findings into Eq. (6.5) renders

$$(6.6) \quad \text{dist}(\hat{\mathbf{z}} - \mathcal{M}) \leq \|\hat{\mathbf{v}} - P_{\mathbf{u}}(\hat{\mathbf{v}})\|_2 + \text{dist}(\mathbf{z} - \mathcal{M}).$$

The proof of the main part of Theorem 6.3.1 follows by applying Corollary 6.4.1 and Lemma 6.4.1 to the first and second terms of the right side of Eq. (6.6), respectively, and noting that  $\tilde{c} = c \langle \mathbf{u}, \hat{\mathbf{u}} \rangle \leq c$ .  $\square$

**6.4.4. Completing the Proof.** In practice, we may not be able to compute the exact value of the bound  $\delta$  in Eq. (6.2), since  $\mathcal{R}$  depends on the third derivative of  $f$ . Instead, we define an approximation of this value and use  $R_{\text{bound}}$  (an upper bound on  $|\mathcal{R}|$ ) to bound the difference between our approximation and  $\delta$ .

Let

$$\delta_{\text{approx}} := c \sqrt{\frac{E^*}{2}} + \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2}.$$

This is the result of setting  $\mathcal{R} = 0$  in Eq. (6.2). We have

$$\begin{aligned}
|\delta - \delta_{\text{approx}}| &= \left| \left[ c \sqrt{\frac{E^*}{2}} + \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R} \right)^{1/2} \right] - \left[ c \sqrt{\frac{E^*}{2}} + \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} \right] \right| \\
&= \left| \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R} \right)^{1/2} - \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} \right|
\end{aligned}$$

## 6.5. REMARKS REGARDING IMPLEMENTATION

---

$$= \begin{cases} \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R} \right)^{1/2} - \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} & \text{if } \mathcal{R} \geq 0 \\ \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} - \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R} \right)^{1/2} & \text{if } \mathcal{R} < 0. \end{cases}$$

Further, since we know that  $|\mathcal{R}| \leq R_{\text{bound}}$ , we have that

$$\left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 - R_{\text{bound}} \right)^{1/2} \leq \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + \mathcal{R} \right)^{1/2} \leq \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + R_{\text{bound}} \right)^{1/2},$$

which produces

$$|\delta - \delta_{\text{approx}}| \leq \begin{cases} \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + R_{\text{bound}} \right)^{1/2} - \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} & \text{if } \mathcal{R} \geq 0 \\ \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} - \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 - R_{\text{bound}} \right)^{1/2} & \text{if } \mathcal{R} < 0. \end{cases}$$

Thus setting

$$\zeta_{\text{bound}} := \max \left\{ \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 + R_{\text{bound}} \right)^{1/2} - \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2}, \right. \\ \left. \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 \right)^{1/2} - \left( \frac{c^4}{4} \sum_{i=2}^m (\kappa^{(i)})^2 - R_{\text{bound}} \right)^{1/2} \right\}$$

implies that  $|\delta - \delta_{\text{approx}}| \leq \zeta_{\text{bound}}$ . This completes the proof of Theorem 6.3.1.  $\square$

## 6.5. Remarks Regarding Implementation

Though Theorem 6.3.1, in particular Lemma 6.4.1, heavily involves a description of the manifold  $\mathcal{M}$  via a function  $\mathbf{f}$ , we stress that  $\mathbf{f}$  does not need to be found explicitly. The properties of  $\mathbf{f}$  stated in Theorem 6.3.1 can be estimated and/or ascertained directly from  $\mathcal{M}$ . In particular, the principal curvatures can be estimated from points on the manifold, and the required continuity and differentiability conditions of  $\mathbf{f}$  at  $\mathbf{x}_0$  depend directly on the smoothness of  $\mathcal{M}$ .

It is important to note that computing the modified Kaslovsky and Meyer tangent bound in Eq. (6.1) may be somewhat difficult in practice due to ascertaining the values of  $n$ ,  $r$ ,  $K$ ,  $\sigma$ ,  $\xi$ , and  $\xi_\lambda$ . However, the authors give some good suggestions for determining these parameters. For

## 6.5. REMARKS REGARDING IMPLEMENTATION

---

example, see their paper [57] for references for estimating the curvatures contained in  $K$  and the noise  $\sigma$ . The radius parameter  $r$ , which denotes the maximum distance between the local origin and the neighbors in local PCA *as measured along the tangent plane*, cannot be computed in practice (because the true tangent plane is unknown). However, Kaslovsky and Meyer [57] suggest that this parameter be estimated using

$$\hat{r}(s) := \sqrt{\frac{1}{2\gamma}(-1 + \sqrt{4\gamma(s - \sigma^2 m)})} \approx r,$$

where  $s$  bounds the squared distance between the neighboring samples and the local origin (e.g.,  $s := \|\mathbf{x}_0 - \mathbf{x}_{n+1}\|_2^2$  using Algorithm 4 in LPCA-SRC) and

$$\gamma := \frac{1}{10}\|\boldsymbol{\kappa}\|_2^2$$

in the case  $d = 1$ . Additionally, the probability constants  $\xi$  and  $\xi_\lambda$  should be chosen to satisfy the constraints of Kaslovsky and Meyer’s (modified) tangent bound and so that the result holds with sufficiently high probability. Lastly, the tangent bound itself provides a means of estimating the number of samples/neighbors parameter  $n$ : Kaslovsky and Meyer state that  $n$  should be chosen to (empirically) minimize the tangent bound [57]. Thus, under generous sampling conditions, doing so could potentially be an alternative to setting  $n$  using cross-validation in LPCA-SRC.

To address the issue of obtaining a clean local origin  $\mathbf{x}_0$ , Kaslovsky and Meyer proposed an efficient algorithm that “tracks” the center of mass of the neighbors of a noisy sample over increasingly large neighborhood radii. The geometric information obtained during this procedure is used to compute a point  $\hat{\mathbf{x}}_0 \approx \mathbf{x}_0$ . See Section 5.1 of their paper [57] for more details. In the context of LPCA-SRC, this algorithm could be used to estimate a clean point from each training sample. The tangent vectors would then be centered at these points instead of the original training samples. We expect that this method would also improve the classification accuracy of LPCA-SRC, though it would increase its computational expense.

Despite these proposed solutions, there is still reason to doubt the practical functionality of Theorem 6.3.1. The information we require from  $\mathcal{M}$  may be impractical to obtain in real-world classification problems for which the true class manifolds are unknown. Though the conditions of Lemma 6.4.1 are satisfied for sufficiently smooth  $\mathcal{M}$  and sufficiently small  $c$ , the precise conditions

## 6.5. REMARKS REGARDING IMPLEMENTATION

---

may be difficult to check when our knowledge of  $\mathcal{M}$  is based entirely on a limited number of sampled points. Further, this problem is exacerbated when the manifold is sparsely-sampled. We admit that this is a substantial shortcoming in the usefulness of Theorem 6.3.1. As a related issue, Theorem 6.3.1 may not hold under the sparse sampling conditions common in classification problems such as face recognition, due to its reliance on Kaslovsky and Meyer’s tangent bound, which requires generous sampling [57]. Another weakness, of course, is that Theorem 6.3.1 is only stated for the case  $d = 1$ . The above analysis is merely a first step towards quantifying the usefulness of the tangent vectors in LPCA-SRC, and more work is needed.

## CHAPTER 7

# Other Local PCA Modifications

In this chapter, we use local PCA to modify two other kinds of representation-based algorithms, analogous to our modification of SRC in LPCA-SRC as presented in Chapter 4. Namely, we consider CRC-RLS [114] (CRC with  $\ell^2$ -regularization) and *classification via structured sparsity* [41] (a method that seeks a *block-sparse* representation of the test sample over the set of training samples). Note that block-sparsity is also referred to as *group-sparsity* in the literature. We compare the original and modified versions of these algorithms with SRC, LPCA-SRC, and SRC<sub>pruned</sub>.

### 7.1. Modifying CRC-RLS

We discussed in Chapter 3 the computational efficiency of using  $\ell^2$ -regularization (instead of  $\ell^1$ -regularization) in the CRC framework. The key to this efficiency is that the matrix  $P := (X_{\text{tr}}^T X_{\text{tr}} + \lambda I)^{-1} X_{\text{tr}}^T$  may be computed offline, and so for any number of test samples, solving Eq. (3.9) with  $p = q = 2$  requires only the matrix-vector multiplication  $\boldsymbol{\alpha}^* = P\mathbf{y}$ .

Since the pruning of the dictionary in LPCA-SRC is computed online, we cannot obtain the same computational speedup in LPCA-SRC by simply replacing the  $\ell^1$ -norm in Eq. (4.2) with the (squared)  $\ell^2$ -norm. In fact, as our experimental results will demonstrate, the matrix multiplication required to compute  $(D_{\mathbf{y}}^T D_{\mathbf{y}} + \lambda I)^{-1} D_{\mathbf{y}}^T$ , even when QR-factorization is used, can be significantly slower than using HOMOTOPY to solve the original  $\ell^1$ -minimization problem, especially when the dictionary is large. However, our primary focus in the following experiments will be on how this non-sparsifying method of regularization (i.e.,  $\ell^2$ -minimization) affects classification accuracy.

**7.1.1. Methods Compared.** We compare SRC, LPCA-SRC, SRC<sub>pruned</sub>, and CRC-RLS with the following modified methods:

- *LPCA-CRC*: LPCA-SRC, but with Eq. (4.2) replaced with

$$(7.1) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\mathbf{y}}}} \left\{ \|\mathbf{y} - D_{\mathbf{y}} \boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_2^2 \right\},$$

## 7.1. MODIFYING CRC-RLS

---

- $CRC_{pruned}$ : SRC<sub>pruned</sub> with the analogous modification,
- $CRC_{unnorm.}$ : This modification of CRC solves

$$(7.2) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{tr}}} \left\{ \|\mathbf{y} - X_{tr}\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_2^2 \right\}$$

and then classifies  $\mathbf{y}$  according to

$$(7.3) \quad \text{class\_label}(\mathbf{y}) = \arg \min_{1 \leq l \leq L} \|\mathbf{y} - X_{tr}\delta_l(\boldsymbol{\alpha}^*)\|_2.$$

Note that  $CRC_{unnorm.}$  is SRC with the  $\ell^1$ -norm replaced with the squared  $\ell^2$ -norm (and without the coefficient of  $1/2$  on the approximation error). Alternatively, it is exactly CRC-RLS without the normalization factor in Eq. (3.10): “unnorm.” abbreviates “unnormalized.”

**7.1.2. Experimental Results.** We compared the above methods on the sinusoidal waveform synthetic database described in Section 5.3 (for a range of training class sizes and noise levels) and on AR-1, the version of the AR database that does not contain images with occlusion. The results were computed over 100 trials for the synthetic database and 10 trials for AR-1. The parameter  $\lambda$  was set to 0.001 in the original algorithms and set using cross-validation in the algorithms LPCA-CRC, CRC-RLS,  $CRC_{unnorm.}$ , and  $CRC_{pruned}$ .

Figure 7.1 shows the accuracy results on the synthetic database as the number of training samples in each class was varied. As we can see, none of the  $\ell^2$ -based algorithms were competitive for small values of  $N_0$ , though both LPCA-CRC and  $CRC_{pruned}$  improved significantly as  $N_0$  increased. The dictionary pruning steps in these methods can be interpreted as an enforcement of (some degree of) sparsity, in the sense that the eliminated training samples are viewed as samples with coefficients equal to 0 in the approximation of  $\mathbf{y}$  over the full dictionary before pruning. This additional aspect of sparsity (for all but  $N_0 = 5$  at which no dictionary pruning occurred) gave LPCA-CRC and  $CRC_{pruned}$  an advantage over the denser  $\ell^2$ -based methods. Even at  $N_0 = 75$ , however, the accuracies of LPCA-CRC and  $CRC_{pruned}$  were still exceeded by those of all three  $\ell^1$ -based algorithms.

## 7.1. MODIFYING CRC-RLS

---

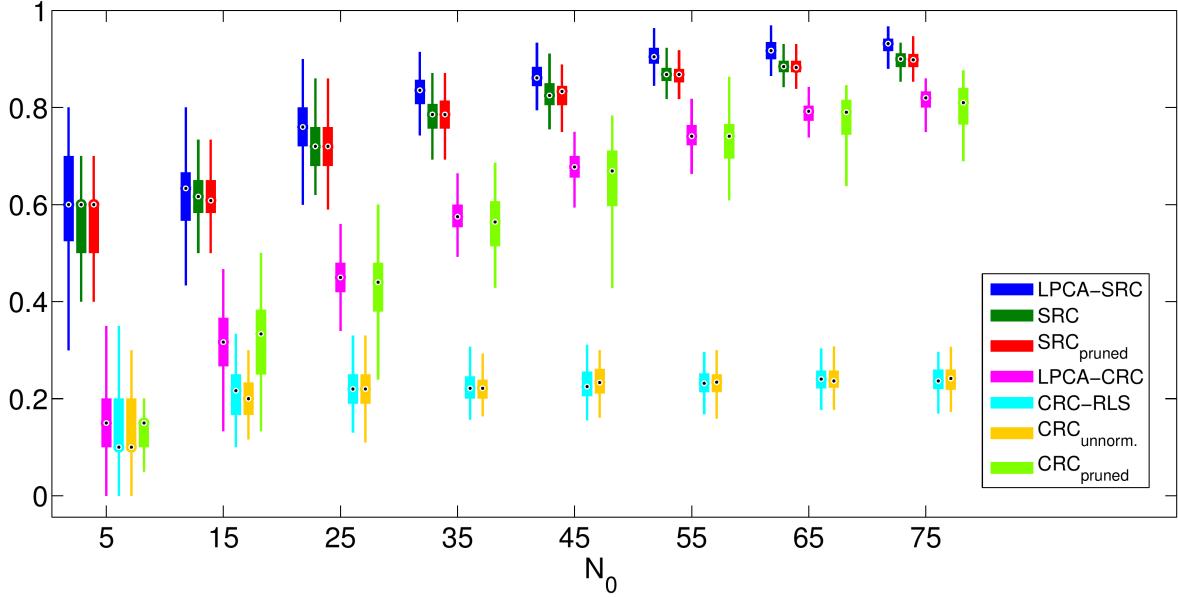


FIGURE 7.1. Box plot of the average classification accuracy (over 100 trials) of the original and  $\ell^2$ -regularized versions of algorithms on the synthetic database with varying training class size  $N_0$ . We fixed  $\eta = 0.001$ .

Though it may be difficult to see in Figure 7.1, LPCA-CRC generally outperformed CRC<sub>pruned</sub>, with the advantage of LPCA-CRC over CRC<sub>pruned</sub> decreasing slightly as  $N_0$  increased, as in the case of the original algorithms LPCA-SRC and SRC<sub>pruned</sub>. The tangent vectors in LPCA-CRC also helped to stabilize classification accuracy, as LPCA-CRC typically had lower variance than CRC<sub>pruned</sub> for sufficiently large  $N_0$ . Lastly, the nearly identical behavior of CRC-RLS and CRC<sub>umnorm.</sub> suggests that the normalization factor in the classification stage of CRC-RLS makes very little difference on this database, though as both methods had classification accuracy comparable to random guessing, this is far from conclusive.

Figure 7.2 shows the results as the noise level  $\eta$  was varied on the synthetic database. Again, though LPCA-CRC and CRC<sub>pruned</sub> performed significantly better than the  $\ell^2$ -based algorithms without dictionary pruning, they were still clearly outperformed by the methods that use  $\ell^1$ -minimization. Though their classification accuracy was only slightly affected by the increase in noise for the displayed values of  $\eta$ , additional experiments confirmed that for  $\eta > 0.5$ , LPCA-CRC and CRC<sub>pruned</sub> continue to be outperformed by the  $\ell^1$ -based methods until all algorithms converge at the point that class structure is entirely lost. Again, the tangent vectors in LPCA-CRC lent it

## 7.1. MODIFYING CRC-RLS

---

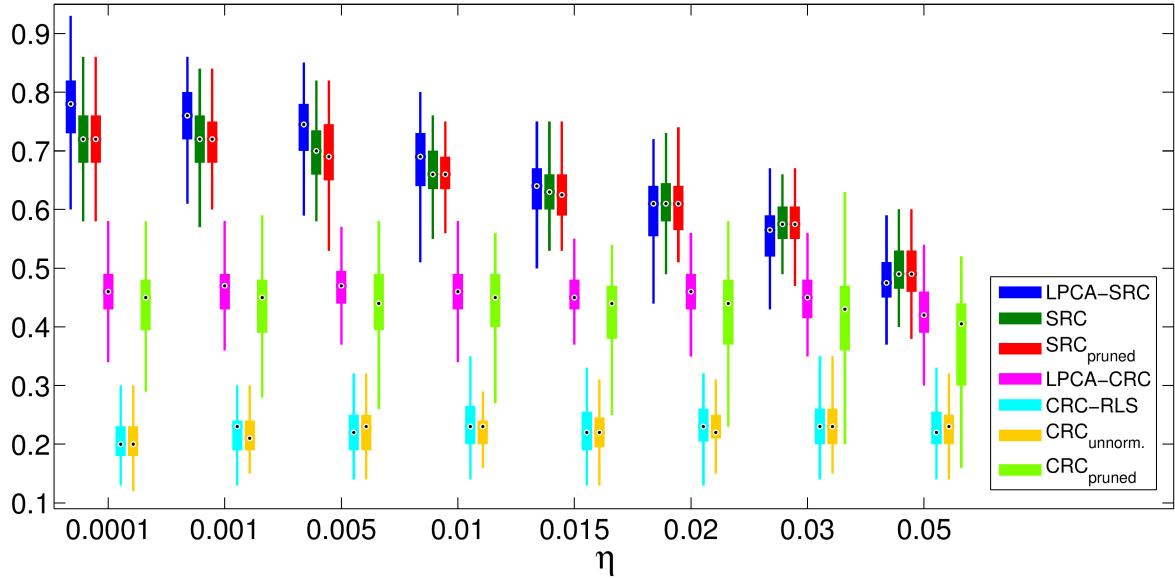


FIGURE 7.2. Box plot of the average classification accuracy (over 100 trials) of the original and  $\ell^2$ -regularized versions of algorithms on the synthetic database with varying noise level  $\eta$ . We fixed  $N_0 = 25$ .

an advantage over  $\text{CRC}_{\text{pruned}}$  (by an average of 3.5%), and the normalization factor in  $\text{CRC-RLS}$  had no discernible effect.

The runtime of each algorithm on the synthetic database is displayed in Table 7.1 for every other value of  $N_0$  and  $\eta$ . As we can see, the  $\ell^2$ -methods generally offered a significant speedup over the methods that use  $\ell^1$ -regularization. The exception, however, was  $\text{LPCA-CRC}$ , which ran in about the same amount of time as  $\text{LPCA-SRC}$ . We also observe that the runtimes of the  $\ell^1$ -methods became larger as  $\eta$  increased. This is the result of a small increase in both the number of iterations and the runtime of each iteration in the  $\ell^1$ -minimization algorithm  $\text{HOMOTOPY}$ , due to the additional noise requiring (relatively) denser coefficient vectors. (The coefficient vectors in each  $\ell^1$ -method had on average 2-3 nonzero coefficients when  $\eta = 0.001$ ; this increased to about 5 nonzero coefficients when  $\eta = 0.05$ .) More details regarding the efficiency of  $\text{HOMOTOPY}$  as it relates to the sparsity of its solution can be found in the paper by Donoho and Tsaig [33]. Note that the closed-form solutions for the coefficient vectors used in the  $\ell^2$ -regularized methods (once their dictionaries were constructed) prevented their runtimes from being similarly affected by the increase in noise.

## 7.1. MODIFYING CRC-RLS

---

Algorithm	Varying $N_0$ ( $\eta = 0.001$ )				Varying $\eta$ ( $N_0 = 25$ )			
	$N_0 = 5$	$N_0 = 25$	$N_0 = 45$	$N_0 = 65$	$\eta = 0.001$	$\eta = 0.01$	$\eta = 0.02$	$\eta = 0.05$
LPCA-SRC	0.011	0.069	0.115	0.137	0.074	0.082	0.096	0.112
SRC	0.004	0.040	0.105	0.126	0.042	0.044	0.052	0.064
SRC <sub>pruned</sub>	0.007	0.054	0.130	0.167	0.056	0.063	0.067	0.077
LPCA-CRC	0.015	0.079	0.100	0.139	0.072	0.101	0.074	0.100
CRC-RLS	0.001	0.005	0.010	0.017	0.005	0.005	0.005	0.005
CRC <sub>unnorm.</sub>	0.001	0.004	0.009	0.015	0.004	0.004	0.004	0.004
CRC <sub>pruned</sub>	0.006	0.045	0.064	0.085	0.042	0.038	0.040	0.047

TABLE 7.1. Average runtime (in seconds) of the original and  $\ell^2$ -regularized versions of algorithms over 100 trials on the synthetic database with varying training class size  $N_0$  and varying noise level  $\eta$ .

We now consider the performance of the various algorithms on the face database AR-1. Table 7.2 contains the average accuracy, standard error, and runtime over 10 trials. We notice that all the  $\ell^2$ -based algorithms became competitive as the dimension increased, performing comparably to their  $\ell^1$ -based counterparts for  $m_{\text{PCA}} = 120$ . Despite their fast runtimes, CRC-RLS and CRC<sub>unnorm.</sub> had significantly worse performance than the other methods for  $m_{\text{PCA}} \in \{30, 56\}$ ; additionally, the normalization factor in CRC-RLS was, if anything, counterproductive. This is in contrast to the CRC authors' experimental finding that the normalization factor slightly improved classification accuracy [114]. It could be the case that this is true at higher feature dimensions or on other databases.<sup>1</sup> On the other hand, CRC<sub>pruned</sub> had comparable accuracy to SRC and SRC<sub>pruned</sub> at all feature dimensions, though it offered very little advantage in runtime. Lastly, as we mentioned at the beginning of this section, LPCA-CRC was *significantly* slower than LPCA-SRC for all tried values of  $m_{\text{PCA}}$ . Though it had accuracy comparable to or exceeding that of SRC and SRC<sub>pruned</sub> at all three feature dimensions, it was still outperformed by the original LPCA-SRC for  $m_{\text{PCA}} \in \{30, 56\}$ . Thus there appear to be only disadvantages—in terms of both classification accuracy and computational efficiency—in using LPCA-CRC instead of LPCA-SRC.

Lastly, we consider the difference in sparsity between the  $\ell^1$ -minimized and  $\ell^2$ -minimized coefficient vectors. Table 7.3 displays the average number of nontrivial coefficients (i.e., coefficients with magnitude at least 0.001) in the solution vector  $\alpha^*$  in the compared algorithms. We normalized

<sup>1</sup>The CRC authors ran experiments on AR-1 with  $m_{\text{PCA}} \in \{54, 120\}$  (as well as  $m_{\text{PCA}} = 300$ ); however, their images were cropped to a different size than we used, and they selected their training set differently than we did [114]. They also did not report the classification accuracy of CRC-RLS without the normalization factor, and so it is difficult to precisely explain our contrasting findings.

## 7.1. MODIFYING CRC-RLS

---

Algorithm	$m_{\text{PCA}} = 30$			$m_{\text{PCA}} = 56$			$m_{\text{PCA}} = 120$		
	Acc	SE	t	Acc	SE	t	Acc	SE	t
LPCA-SRC	<b>0.8663</b>	4.1	7.25	<b>0.9544</b>	2.3	12.50	0.9711	1.7	19.07
SRC	0.8273	4.2	6.11	0.9357	2.6	8.87	0.9631	1.6	13.57
$\text{SRC}_{\text{pruned}}$	0.8277	4.8	3.76	0.9353	3.8	5.10	0.9651	1.8	6.90
LPCA-CRC	0.8339	6.9	121.09	0.9409	1.5	385.81	<b>0.9729</b>	2.1	269.69
CRC-RLS	0.7516	4.9	0.64	0.9080	3.6	0.67	0.9611	2.0	0.74
$\text{CRC}_{\text{unnorm.}}$	0.7584	5.2	0.54	0.9070	3.1	0.58	0.9650	2.5	0.63
$\text{CRC}_{\text{pruned}}$	0.8357	4.4	5.51	0.9366	2.1	5.45	0.9659	2.1	5.57

TABLE 7.2. Average accuracy, standard error ( $\times 10^{-3}$ ), and runtime (in seconds) of the original and  $\ell^2$ -regularized versions of algorithms over 10 trials on AR-1.

this number by the length of the (respective) dictionary. Since the coefficient vectors are the same in CRC-RLS and  $\text{CRC}_{\text{unnorm.}}$ , we only show the results for the latter algorithm.

Algorithm	Synthetic ( $\eta = 0.001$ )				Synthetic ( $N_0 = 25$ )				AR-1		
	$N_0 = 5$	25	45	65	$\eta = 0.001$	0.01	0.02	0.05	$m_{\text{PCA}} = 30$	56	120
LPCA-SRC	0.15	0.12	0.21	0.19	0.13	0.06	0.08	0.09	0.22	0.53	0.65
LPCA-CRC	0.98	0.98	0.99	0.99	0.99	0.97	0.98	0.98	0.78	0.74	0.82
SRC	0.12	0.13	0.07	0.04	0.19	0.04	0.04	0.06	0.19	0.40	0.82
$\text{CRC}_{\text{unnorm.}}$	0.99	0.97	0.94	0.92	0.97	0.98	0.98	0.99	0.87	0.90	0.92
$\text{SRC}_{\text{pruned}}$	0.13	0.14	0.08	0.02	0.17	0.05	0.07	0.13	0.47	0.84	0.96
$\text{CRC}_{\text{pruned}}$	0.99	0.99	0.99	1.00	0.99	0.99	0.99	0.99	0.95	0.96	0.97

TABLE 7.3. Average number of nontrivial coefficients (divided by the total number of coefficients) in the original and  $\ell^2$ -regularized versions of the algorithms on the synthetic database (100 trials) and AR-1 (10 trials).

It should come as no surprise that on all versions of the synthetic database, the  $\ell^1$ -minimized coefficient vectors were *significantly* sparser than their  $\ell^2$ -based counterparts. The results on AR-1, however, are intriguing: Though the coefficient vectors of the original methods LPCA-SRC, SRC, and  $\text{SRC}_{\text{pruned}}$  were relatively sparse at  $m_{\text{PCA}} = 30$  (for example, the methods LPCA-SRC and SRC used only a quarter of the nontrivial coefficients required by LPCA-CRC and  $\text{CRC}_{\text{unnorm.}}$ , respectively), they became denser as the feature dimension increased. At  $m_{\text{PCA}} = 120$ , for instance,  $\text{SRC}_{\text{pruned}}$  and  $\text{CRC}_{\text{pruned}}$  both had nearly maximally-dense coefficient vectors. The reason for this increasing number of nontrivial coefficients in the  $\ell^1$ -based algorithms is likely due to the redundancy of the dictionary decreasing (i.e.,  $m_{\text{PCA}}$  becoming closer to  $N$ ). This helps to explain why the  $\ell^2$ -regularized methods performed comparably to the original algorithms at  $m_{\text{PCA}} = 120$ .

REMARK 7.1.1. We could have solved the optimization problem in all compared algorithms using Friedman et al.’s *glmnet* package [43, 86], a coordinate-descent approach to solving generalized linear models<sup>2</sup> with either  $\ell^1$ -regularization,  $\ell^2$ -regularization, or a combination of both (this is called elastic-net). In the case of LPCA-CRC especially, this approach is much more efficient than the closed-form matrix multiplication discussed above. However, it is recommended that the columns of the input dictionary in *glmnet* be centered so that their means are 0.<sup>3</sup> Since the mean of each sample (generally) contains class information, this pre-processing significantly decreased the classification accuracy of LPCA-SRC, SRC, and SRC<sub>pruned</sub> compared to our experiments in Chapter 5 using HOMOTOPY.

An additional deterrent in using *glmnet* to solve the optimization problems in representation-based classification algorithms involves the specification of the regularization parameter  $\lambda$ . Designed primarily for regression applications, *glmnet* does not generally work well with values of  $\lambda$  too close to 0, as the resulting solution would generally correspond to overfitting. In the case that we want  $\mathbf{y}$  to be a nearly-exact linear combination of the dictionary elements (i.e.,  $\lambda \approx 0$ ), proper use of *glmnet* will often return a solution corresponding to an undesirably large value of  $\lambda$ .

## 7.2. Modifying Structured Sparsity/Block-Sparse Methods

In SRC, we assume that the test sample can be well-approximated by its same class training samples, with the goal being that nonzero coefficients in the  $\ell^1$ -minimized approximation of  $\mathbf{y}$  occur primarily at training samples in that class. However, as the desired coefficient vector has the form

$$\boldsymbol{\alpha} = [0, \dots, 0, \alpha_1^{(l)}, \dots, \alpha_{N_l}^{(l)}, 0, \dots, 0]^T \in \mathbb{R}^{N_{\text{tr}}}$$

for  $\mathbf{y}$  in class  $l$ , seeking the sparsest vector is not necessarily the most direct approach, especially if the number of classes  $L$  is small. Instead, one could seek a representation of the test sample that constrains the number of *classes* that contain training samples with nonzero coefficients. In other words, given that the columns of  $X_{\text{tr}}$  are grouped by class so that we can write  $X_{\text{tr}} = [X^{(1)}, \dots, X^{(L)}]$ , then if we partition  $\boldsymbol{\alpha}$  into the “blocks”  $\boldsymbol{\alpha} = [(\boldsymbol{\alpha}^{(1)})^T, \dots, (\boldsymbol{\alpha}^{(L)})^T]^T$  where  $\boldsymbol{\alpha}^{(l)} =$

<sup>2</sup>The least squares terms in Eq. (3.7) and Eq. (7.2), for example, make these *linear* models.

<sup>3</sup>Centering is necessary in our application because it ensures that the *regression intercept* in the *glmnet* optimization problem is equal to 0, as is required for representation-based classification. See Hastie et al.’s book [49] for more details.

## 7.2. MODIFYING STRUCTURED SPARSITY/BLOCK-SPARSE METHODS

---

$[\alpha_1^{(l)}, \dots, \alpha_{N_l}^{(l)}]^T \in \mathbb{R}^{N_l}$ ,  $1 \leq l \leq L$ , then we could aim to minimize the number of blocks  $\boldsymbol{\alpha}^{(l)}$  that contain nonzero coefficients. Of course, the ideal solution has nonzero coefficients in exactly one block, and it is assumed that this block corresponds to the correct class.

In one variant of the *structured sparse representation* (SSR) classification algorithm of Elhamifar and Vidal [40], the optimization problem behind this concept is written as

$$(7.4) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \sum_{l=1}^L \mathbb{I}(\|\boldsymbol{\alpha}^{(l)}\|_2 > 0) \text{ subject to } \mathbf{y} = X_{\text{tr}}\boldsymbol{\alpha}.$$

However, as in the case of finding the sparsest coefficient vector (via  $\ell^0$ -“norm” minimization), solving Eq. (7.4) is intractable. Fortunately, the relaxation

$$(7.5) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \sum_{l=1}^L \|\boldsymbol{\alpha}^{(l)}\|_2 \text{ subject to } \mathbf{y} = X_{\text{tr}}\boldsymbol{\alpha}$$

has been proven to recover the solution to Eq. (7.4) under certain conditions [37, 41]. Note that the objective function is a special case of the so-called *mixed  $\ell^p/\ell^q$ -norm*:

$$\left( \sum_{l=1}^L \|\boldsymbol{\alpha}^{(l)}\|_p^q \right)^{1/q},$$

with  $p = 2$  and  $q = 1$ . Much as minimizing the  $\ell^1$ -norm can be used to determine the sparsest solution in certain situations, minimizing the mixed  $\ell^2/\ell^1$ -norm can be used to retrieve what we call the *block-sparsest solution*, i.e., the solution with nonzeros in the minimal number of blocks [37]. This is referred to as *block-sparse recovery* (BSR). Classification in SSR is then determined using the solution  $\boldsymbol{\alpha}^*$  of Eq. (7.5) in Eq. (3.6), as in SRC.

When noise makes the exact representation  $\mathbf{y} = X_{\text{tr}}\boldsymbol{\alpha}$  impractical, we can instead solve

$$(7.6) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \|\mathbf{y} - X_{\text{tr}}\boldsymbol{\alpha}\|_2^2 + \lambda \sum_{l=1}^L \|\boldsymbol{\alpha}^{(l)}\|_2.$$

In this format, BSR is easily identifiable as the *group-LASSO* (*Least Absolute Shrinkage and Selection Operator*), which was originally proposed by Yuan and Lin to solve regression problems [112].

## 7.2. MODIFYING STRUCTURED SPARSITY/BLOCK-SPARSE METHODS

---

On the other hand, Eq. (7.6) is equivalent to

$$(7.7) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \sum_{l=1}^L \|\boldsymbol{\alpha}^{(l)}\|_2 \text{ subject to } \|\mathbf{y} - X_{\text{tr}}\boldsymbol{\alpha}\|_2 \leq \epsilon,$$

the formulation used for robust block-sparse signal recovery by Eldar and Mishali [37].

**7.2.1. Methods Compared.** We compared SRC, LPCA-SRC, and SRC<sub>pruned</sub> with this version of SSR, formulating the optimization problem as in Eq. (7.7). For fair comparison, we computed new results for the first three algorithms with their respective  $\ell^1$ -minimization problems formulated as

$$(7.8) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \|\boldsymbol{\alpha}\|_1 \text{ subject to } \|\mathbf{y} - D\boldsymbol{\alpha}\|_2 \leq \epsilon,$$

where  $D \in \mathbb{R}^{m \times N}$  denotes the algorithm's dictionary. Note that this formulation is equivalent to each method's previously presented optimization problem under an appropriate correspondence between the parameters  $\epsilon$  and  $\lambda$ ; however, this correspondence cannot be determined a priori.

We also considered the relative classification performance of the modified algorithms LPCA-SSR, wherein Eq. (4.2) is replaced with

$$(7.9) \quad \boldsymbol{\alpha}^* := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \sum_{l=1}^L \|\boldsymbol{\alpha}^{(l)}\|_2 \text{ subject to } \|\mathbf{y} - D_y \boldsymbol{\alpha}\|_2 \leq \epsilon,$$

and SSR<sub>pruned</sub>, which is SRC<sub>pruned</sub> with a similar revision.

To solve the optimization problems in all six methods, we used the sparse reconstruction toolbox SPGL1 [95, 96], due to its easily-implementable methods for solving both Eq. (7.7) and Eq. (7.8). The parameter  $\epsilon$  was set using cross-validation in all classification algorithms.

**7.2.2. Experimental Results.** As in the previous section, we compared the above algorithms on the sinusoidal waveform synthetic database from Chapter 5 and the face database AR-1.

Figure 7.3 shows the results as the training class size  $N_0$  was varied on the synthetic database. The immediate observation is that SSR performed significantly worse than SRC and the other  $\ell^1$ -based methods. Though the theoretical motivation behind SSR and that of CRC-RLS discussed in the previous section are completely different, their reasons for failure on this database are essentially the same: since each training class can generally provide a decent approximation of the test sample,

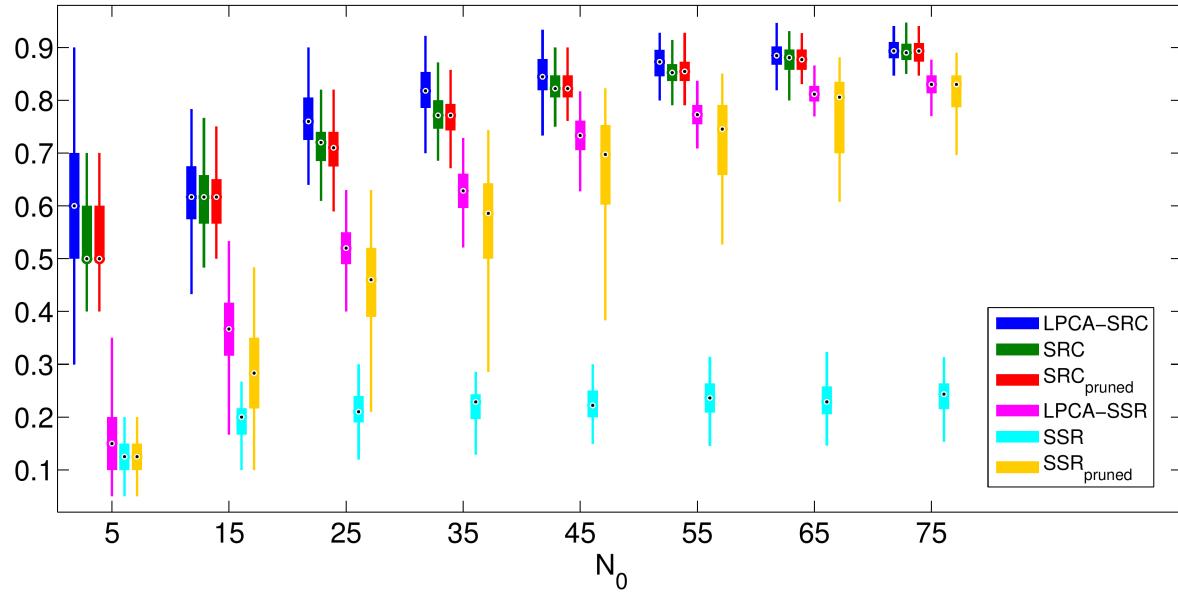


FIGURE 7.3. Box plot of the average classification accuracy (over 100 trials) of the original and SSR versions of algorithms on the synthetic database with varying training class size  $N_0$ . We fixed  $\eta = 0.001$ .

the classes must be differentiated *locally*, and these algorithms do not have the machinery to do this. The block-sparse prior behind SSR only considers each class as a whole, and so it cannot identify the class manifold structure in a small neighborhood around  $\mathbf{y}$ . As in the case of the  $\ell^2$ -regularized algorithms, the SSR methods that utilize dictionary pruning, i.e., LPCA-SSR and SSR<sub>pruned</sub>, performed notably better than the other block-sparse methods. LPCA-SSR performed better than SSR<sub>pruned</sub> (in terms of both classification accuracy and variance), though as we would expect, its advantage decreased as the sampling density increased. Further, neither method could match the performance of the algorithms that use  $\ell^1$ -minimization, presumably because the pruned dictionary still allowed for a good approximation of the test sample in terms of each class separately. Note that if the dictionary pruning procedure were stricter or if the frequency of the sinusoidal waves was smaller (see the database description in Section 5.3), this might not be the case, and we would expect these methods to match or exceed their  $\ell^1$ -based counterparts in terms of accuracy.

Figure 7.4 shows the results as the noise level  $\eta$  was varied on the synthetic database. Our observations are nearly-identical to those made of Figure 7.3, and so we do not restate them.

## 7.2. MODIFYING STRUCTURED SPARSITY/BLOCK-SPARSE METHODS

---

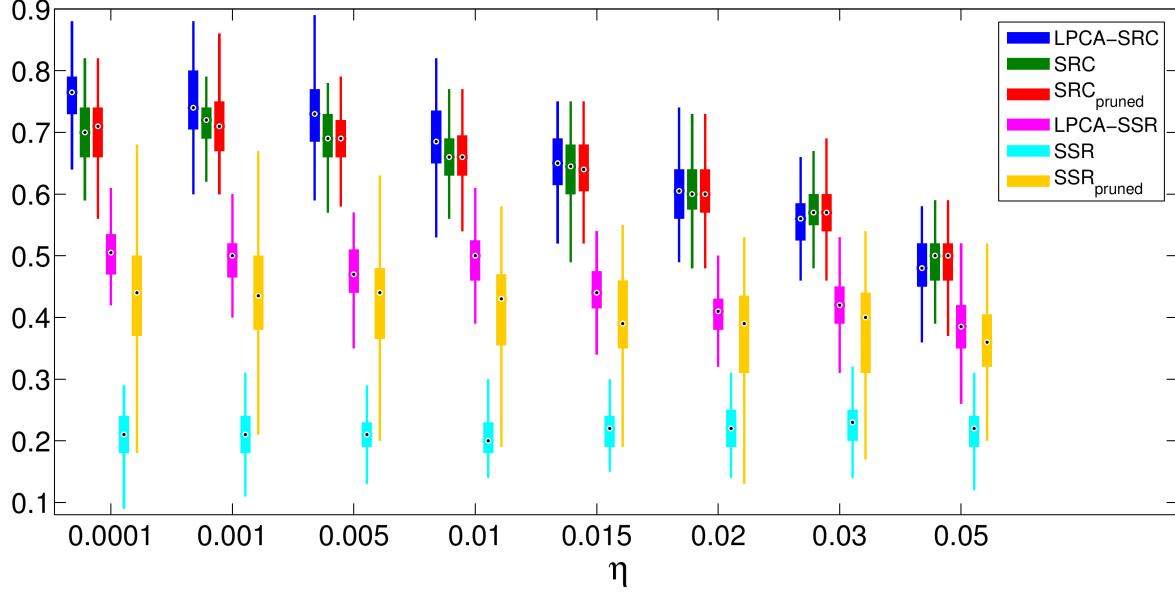


FIGURE 7.4. Box plot of the average classification accuracy (over 100 trials) of the original and SSR versions of algorithms on the synthetic database with varying noise level  $\eta$ . We fixed  $N_0 = 25$ .

To assess computational efficiency, the runtime of each algorithm on the synthetic database is displayed in Table 7.4 for every other value of  $N_0$  and  $\eta$ . It is clear that, using SPGL1 on this database, the runtime to solve the BSR optimization problem scales poorly with  $N_{\text{tr}}$ , relative to the runtime of determining the  $\ell^1$ -minimization solution.

Algorithm	Varying $N_0$ ( $\eta = 0.001$ )				Varying $\eta$ ( $N_0 = 25$ )			
	$N_0 = 5$	$N_0 = 25$	$N_0 = 45$	$N_0 = 65$	$\eta = 0.001$	$\eta = 0.01$	$\eta = 0.02$	$\eta = 0.05$
LPCA-SRC	0.27	1.72	2.67	2.83	1.86	0.61	0.59	0.83
SRC	0.16	2.89	3.99	4.33	2.34	0.64	0.57	0.58
SRC <sub>pruned</sub>	0.14	2.23	3.81	2.91	2.79	0.67	0.57	0.57
LPCA-SSR	0.55	5.89	11.82	18.35	5.94	1.93	3.03	3.13
SSR	0.17	2.35	6.02	9.49	2.15	5.36	5.64	3.79
SSR <sub>pruned</sub>	0.25	2.01	7.76	12.90	3.33	4.06	4.56	3.38

TABLE 7.4. Average runtime (in seconds) of the original and SSR versions of algorithms over 100 trials on the synthetic database with varying training class size  $N_0$  and varying noise level  $\eta$ .

We next consider the performance of the various algorithms on the face database AR-1. Table 7.5 contains the average accuracy, standard error, and runtime over 10 trials. We continue to see

## 7.2. MODIFYING STRUCTURED SPARSITY/BLOCK-SPARSE METHODS

---

the  $\ell^1$ -based methods obtain higher classification accuracy than the methods that aim for block-sparsity. The latter algorithms improved as the dimension increased, though the  $\ell^1$ -based methods still maintained the advantage at  $m_{\text{PCA}} = 120$ . This is discussed further throughout the rest of this chapter. With regard to runtime, it appears that the SSR methods scale much better with the feature dimension  $m_{\text{PCA}}$  than the training set size  $N_{\text{tr}}$  (for which poor scaling was observed on the synthetic database): the modified algorithms ran in roughly the same amount of time as the  $\ell^1$ -based algorithms for all three values of  $m_{\text{PCA}}$ .

We also see, interestingly, that using SPGL1 in the LPCA-SRC and SRC algorithms produced slightly different results than using HOMOTOPY: LPCA-SRC had a slightly larger lead over SRC at  $m_{\text{PCA}} = 30$  with SPGL1 (5% instead of 4%); however, SRC actually outperformed LPCA-SRC at  $m_{\text{PCA}} = 120$  by 2%, instead of a 1% difference in the opposite direction using HOMOTOPY. This is likely the result of setting  $\epsilon$  using cross-validation in SPGL1, whereas  $\lambda$  was fixed at 0.001 using HOMOTOPY on this database.

Algorithm	$m_{\text{PCA}} = 30$			$m_{\text{PCA}} = 56$			$m_{\text{PCA}} = 120$		
	Acc	SE	t	Acc	SE	t	Acc	SE	t
LPCA-SRC	<b>0.8499</b>	12.5	52.52	<b>0.9451</b>	6.8	108.33	0.9457	8.9	149.05
SRC	0.7969	5.6	63.35	0.9244	3.8	120.45	<b>0.9667</b>	2.5	189.97
SRC <sub>pruned</sub>	0.8211	5.1	44.37	0.9320	2.5	77.34	0.9514	3.0	125.99
LPCA-SSR	0.7587	3.9	55.07	0.9020	4.4	102.90	0.9421	7.0	166.91
SSR	0.6516	6.2	60.12	0.8309	4.8	112.92	0.9276	3.2	157.20
SSR <sub>pruned</sub>	0.7436	6.3	48.19	0.8946	9.0	77.08	0.9343	11.0	141.03

TABLE 7.5. Average accuracy, standard error ( $\times 10^{-3}$ ), and runtime (in seconds) of the original and SSR versions of algorithms over 10 trials on AR-1.

Lastly, we consider the question of whether the optimization problem in Eq. (7.7) really recovered a block-sparse solution in the above experiments. Table 7.6 displays the average number of classes that contained training samples with nontrivial coefficients (i.e., coefficients with magnitude at least 0.001) in the solution vector  $\alpha^*$  for each of the six compared algorithms. We normalized this number by the total number of classes in the database.

We observe that the SSR methods were not, in general, more block-sparse than the methods that use  $\ell^1$ -minimization, and they were sometimes less so. This is particularly clear in comparing SRC and SSR: the former algorithm was more block-sparse (sometimes significantly) at large values

## 7.2. MODIFYING STRUCTURED SPARSITY/BLOCK-SPARSE METHODS

---

Algorithm	Synthetic ( $\eta = 0.001$ )				Synthetic ( $N_0 = 25$ )				AR-1		
	$N_0 = 5$	25	45	65	$\eta = 0.001$	0.01	0.02	0.05	$m_{PCA} = 30$	56	120
LPCA-SRC	0.53	0.44	0.44	0.41	0.51	0.50	0.57	0.72	0.38	0.47	0.56
LPCA-SSR	0.69	0.58	0.58	0.47	0.69	0.51	0.52	0.64	0.35	0.50	0.51
SRC	0.58	0.39	0.39	0.32	0.59	0.53	0.56	0.74	0.63	0.91	1.00
SSR	0.47	0.46	0.46	0.65	0.45	0.75	0.86	0.84	0.61	0.95	1.00
SRC <sub>pruned</sub>	0.56	0.39	0.39	0.40	0.52	0.54	0.58	0.76	0.50	0.60	0.61
SSR <sub>pruned</sub>	0.49	0.50	0.50	0.45	0.57	0.62	0.62	0.66	0.45	0.55	0.58

TABLE 7.6. Average number of classes (divided by the total number of classes) with nontrivial coefficients in compared algorithms on the synthetic database (10 trials) and AR-1 (10 trials).

of both  $N_0$  and  $\eta$  on the synthetic database, though both methods had similar block-sparsity on AR-1 and placed nontrivial coefficients at nearly all of the classes for  $m_{PCA} = 56$  and at all of the classes for  $m_{PCA} = 120$  on this data set. The methods SRC<sub>pruned</sub> and SSR<sub>pruned</sub> generally had about the same amount of block-sparsity; however SRC<sub>pruned</sub> was often slightly more block-sparse than SSR<sub>pruned</sub> on the synthetic database, and in contrast, SSR<sub>pruned</sub> was slightly block-sparser on AR-1. LPCA-SRC and LPCA-SSR also reported similar levels of block-sparsity, though the former algorithm was slightly more block-sparse on the synthetic database, except at large values of  $\eta$  for which the opposite was true. These differences were relatively small. The important observation is that we do not see convincing recovery of the block-sparsest solution via Eq. (7.7) in these results.

Further, as the values for AR-1 in Table 7.6 demonstrate, the degree of block-sparsity of each algorithm did not always correspond to its ability to obtain relatively high classification accuracy. Consider LPCA-SRC, for example, which had nontrivial coefficients at many fewer classes than SRC for  $m_{PCA} = 120$ , yet still the latter algorithm had higher classification accuracy at this feature dimension. Though our results rely on the threshold we used to define “nontrivial” coefficient values, this observation is perhaps important to explaining why  $\ell^1$ -methods may sometimes work better than those based on BSR. This explanation, however, requires further investigation into block-sparse recovery theory: how close is the minimal mixed  $\ell^2/\ell^1$ -norm solution in Eq. (7.7) to the ground truth block-sparsest solution on these databases? Like the relationship between

### 7.3. DISCUSSION

---

SRC and the sparsest solution, we suspect that it may be difficult to verify block-sparse recovery conditions in these circumstances.<sup>4</sup>

#### 7.3. Discussion

The synthetic database is notable in the context of representation-based algorithms in the sense that the corresponding system of equations is *very* underdetermined, with the set of training samples from any single class often able to produce a good approximation of the test sample. Geometrically, identifying the correct class comes down to isolating an *extremely* local neighborhood around  $\mathbf{y}$ . Neither  $\ell^2$ -minimization nor block-sparsity provided a means of doing this. As in the case of LSDL-SRC in Chapter 4, the dictionary pruning step in LPCA-CRC,  $\text{CRC}_{\text{pruned}}$ , LPCA-SSR and  $\text{SSR}_{\text{pruned}}$  managed to enforce some degree of locality, which improved classification performance. However, these methods could not match the effectiveness of  $\ell^1$ -minimization, which achieved exceptionally local approximations of the test sample through its relationship to sparsity. Further, in the case of the SSR methods, the results in Table 7.6 suggest that minimization of the mixed  $\ell^2/\ell^1$ -norm was sometimes unable to recover the block-sparsest solution in this very underdetermined case.

On the face database AR-1, the negative correlation between the relative advantage of the  $\ell^1$ -based methods and the feature dimension  $m_{\text{PCA}}$  is intriguing. On this database,  $\ell^1$ -minimization (due to its relationship with sparsity) was critical only when  $m_{\text{PCA}}$  was small. In the case of the  $\ell^2$ -regularized algorithms, this has been (imprecisely) explained in terms of the discriminative information contained in the test sample at large feature dimensions. The CRC authors Zhang et al. argued that  $\ell^2$ -minimization is sufficient in these cases to produce a *naturally* sparse coefficient vector without the heavy machinery of  $\ell^1$ -minimization [115]. Our experimental results in Table 7.3 suggest, however, that the  $\ell^2$ -minimized coefficient vectors are *not* actually sparse, and that the diminishing value of  $\ell^1$ -minimization as the feature dimension increases is due to the decreased redundancy in the dictionary resulting in a loss of sparsity in these algorithms.

With regard to block-sparse methods, we stress that these limited experiments on the synthetic database and AR-1 are insufficient to draw general conclusions about the efficacy of this approach

---

<sup>4</sup>In particular, we conjecture that this is the case in applications such as face recognition in which the classes themselves are highly-correlated. We discuss this correlation further in Chapter 11.

### 7.3. DISCUSSION

---

to classification. The particular class manifold structure of the synthetic database is clearly biased against methods that do not differentiate the training samples *within* each class, and it could be the case that we would see better block-sparse recovery in SSR on AR-1 at higher feature dimensions, under different pre-processing, when there are less classes, or when a different optimization algorithm is used. Again, full understanding of why  $\ell^1$ -methods may perform better than BSR methods, as is the case of AR-1, involves investigation of how well real-world, class-structured data satisfies block-sparse recovery guarantees. This was not discussed in the SSR paper [40].

As a concluding remark, it would be informative to consider the relative performance of classifiers based on combinations of these methods. According to Hastie et al. [49], elastic-net regularization, which penalizes both the  $\ell^1$ -norm and  $\ell^2$ -norm of the coefficient vector, can improve performance in regression applications over  $\ell^1$ -regularization alone. Elastic-net tends to return nonzero coefficients in “blocks” of correlated variables, i.e., selecting either all variables in the block or none of them [49]. Since Hastie et al.’s regression problem is formulated so that variables correspond to columns of the dictionary, this suggests that elastic-net could be used in representation-based classification to impose a (perhaps weak) block-sparse prior. When the blocks of highly-correlated samples correspond to classes, this method could be effective. Note that this is not a novel observation: see, for example, the paper by Timofte and Van Gool [91].

Methods that combine  $\ell^1$ -minimization with BSR have also been proposed, which force the coefficient vector to be both block-sparse and sparse within each block. This approach is argued to, among other things, improve robustness. See, for example, Sprechmann et al.’s *Hierarchical Lasso* (HiLasso) [88], which has been applied to dictionary learning [90] and multiview facial recognition [117]. On the synthetic database especially, we suspect that this approach would be quite effective.

## CHAPTER 8

### Conclusion of Part 1

The main contribution of Part 1 was a modification of SRC called *local principal component analysis SRC*, or “LPCA-SRC.” Through the use of tangent vectors, LPCA-SRC is designed to increase the sampling density of training sets and thus improve class discrimination in SRC on databases with sparsely-sampled and/or nonlinear class manifolds. The LPCA-SRC algorithm computes basis vectors of approximate tangent hyperplanes at the training samples in each class and replaces the dictionary of training samples in SRC with a local dictionary computed from shifted and scaled versions of these basis vectors and their corresponding training samples. Using a synthetic database and three face databases, we showed that LPCA-SRC can regularly achieve higher accuracy than SRC in cases of sparsely-sampled and/or nonlinear class manifolds, low noise, and relatively small PCA feature dimension.

Though the implementation details, particularly the setting of parameters, is somewhat involved in LPCA-SRC, the general concept of filling out class manifolds via tangent vectors is straightforward and easy to illustrate in low dimension. To address the issue of parameter setting, we suggested several viable methods for setting the class manifold dimension estimate  $d$ . It is important to note that in the case of small training sets, e.g., many face recognition problems, there are few options for the number-of-neighbors parameter  $n$ —and consequently for  $d$  by Eq. (4.3)—and so these values can easily be set using cross-validation, as in our experiments. When the training sets are very small (i.e.,  $N_l = 4$  or  $5$ ), one could simply set the parameters  $d$  and  $n$  to their maximum values, i.e.,  $d = n = N_{l_{\min}} - 2$ , per Eq. (4.3). On the other hand, when algorithm efficiency is paramount, simply setting  $d = 1$  may suffice.

One disadvantage of this method is its high computational cost. SRC is already expensive due to its  $\ell^1$ -minimization procedure; in LPCA-SRC, the computation of tangent vectors is added to the algorithm’s workload. The size of the dictionary in LPCA-SRC may be larger or smaller than that of SRC, depending on the LPCA-SRC parameters  $n$  and  $d$  and the effect of the pruning

---

parameter  $r$ . Thus LPCA-SRC can be slower or faster than SRC. As mentioned in Section 5.4.7, in deciding whether or not to use LPCA-SRC instead of SRC, simple computations based on the training data could be used to estimate the size of the pruned dictionary in LPCA-SRC: the user could compute the pruning parameter  $r$ , determine the number of training samples within  $r$  of  $\mathbf{y}$ , and then multiply this number by  $d + 1$ . Knowledge of the dictionary size would give the user an indication of how the runtimes of LPCA-SRC and SRC would compare on the given data set.

Additionally, as we saw on the synthetic database, the usefulness of the tangent vectors in LPCA-SRC decreases as the noise level in the training data increases. This problem could likely be alleviated by using the method proposed by Kaslovsky and Meyer [57] to estimate clean points on the manifolds from noisy samples and then computing the tangent vectors at these points, as discussed in Chapter 6. Under the assumption that an approximately clean point is determined, we stated a theorem bounding (whp) the distance between an LPCA-SRC tangent vector and the corresponding class manifold, giving the user a potential tool to estimate the noise level of the computed tangent vectors when viewed as additional training samples. This information could then be factored into the decision of whether or not to use LPCA-SRC instead of SRC. Note that the case of large training sample noise was the only case for which we saw LPCA-SRC not obtain higher accuracy than SRC. Thus LPCA-SRC should be preferred over SRC in low noise scenarios on either small-scale problems (e.g., the size of ORL) or when achieving a modest (e.g., 1% – 4%) boost in accuracy is worth potentially higher computational cost.

With regards to the *interpretability* (as discussed in Chapter 2) of LPCA-SRC, as for all representation-based classifiers, the position of large (in magnitude) coordinates of the coefficient vector  $\boldsymbol{\alpha}^*$  can be used to identify training samples that are similar to the test sample, in the sense that they contribute substantially to its approximation. Additionally, the dictionary pruning technique in LPCA-SRC can give us some indication of the *separability* of the class manifolds via the average number of classes represented in the dictionary  $D_{\mathbf{y}}$ : if few, then this suggests that the class manifolds are fairly well-separated using Euclidean distance. Finally, visualization of the tangent vectors, as in Figures 5.4c, 5.5c, and 5.6c, can reveal details lost in PCA dimensionality reduction, as well as a small amount of information regarding the structure of each class manifold.

---

To test the generalization ability of our local PCA modification, we similarly used tangent vectors and dictionary pruning to amend CRC-RLS and Elhamifar and Vidal’s SSR classification algorithm based on block sparsity. In both cases, our modification substantially improved performance. However, because the dictionary pruning step in our framework is performed online, LPCA-CRC was much slower than CRC-RLS (and no faster than the original LPCA-SRC), especially on large databases. Our experiments also showed that on the synthetic database and AR-1 for small feature dimensions, methods based on  $\ell^1$ -minimization significantly outperformed their  $\ell^2$ -regularized and block-sparse counterparts. In the case of the latter type of method, it could be that minimization of the mixed  $\ell^2/\ell^1$ -norm, used as an attempt at determining the block-sparest solution, had poor recovery on these databases. However, given the motivation behind the BSR methods (in that they are *specifically* aimed at placing nonzero coefficients at training samples in the correct class), our empirical results should not discount BSR methods in general, especially those that simultaneously utilize  $\ell^1$ -minimization.

Open questions regarding LPCA-SRC include whether or not the aforementioned general trends hold on other databases, in particular, in different classification applications. On the USPS database [54] for example, which contains samples of handwritten digits, LPCA-SRC does not offer an advantage over SRC, partially because significantly more energy is contained in the PCA feature vectors in USPS than those in the above face databases when the PCA dimension is the same.<sup>1</sup> Though LPCA-SRC is designed to increase the classification accuracy of SRC in the case of sparsely-sampled and/or nonlinear class manifolds, this algorithm requires a minimum sampling density (relative to the curvature) in order to obtain good tangent hyperplane estimates, as discussed in Chapter 6. Further, we do not expect it to perform well when the class manifolds are not smooth. On the other hand, if the class manifolds are *too* well-sampled with respect to curvature, tangent vectors are unneeded and LPCA-SRC will not offer an improvement over SRC. Since experimental results have indicated that handwritten-digit manifolds can be represented as lower-dimensional subspaces without losing too much information [19, 38], we suspect that this is why LPCA-SRC does not outperform SRC on the well-sampled USPS database. However, since the class manifold

---

<sup>1</sup>For the curious reader, we note that LPCA-CRC and all three BSR methods were nearly as good (within roughly 2%) of the LPCA-SRC and SRC results on USPS. CRC was slightly worse, though it did offer an improvement over  $\text{CRC}_{\text{unnorm.}}$  and  $\text{CRC}_{\text{pruned}}$ .

---

structure can only be roughly estimated (at best) in most real-world applications, it is difficult to predict the efficacy of LPCA-SRC in general.

Finally, it remains to investigate whether LPCA-SRC can improve upon SRC when different methods of dimensionality reduction besides PCA are used. One could also compare the performance of the per-class decomposition versions (in which test samples are approximated using class-specific dictionaries as discussed in Section 3.1) of the representation-based algorithms used in our experiments.

## **Part 2**

### **Examining Sparsity in Classification**

## CHAPTER 9

### Equivalence Guarantees

In the second part of this dissertation, we turn to the assumption on which SRC is fundamentally based: namely, that  $\ell^1$ -minimization recovers a sparse solution. In this chapter, we review background material regarding so-called  $\ell^1/\ell^0$ -*equivalence* (referring to conditions under which minimization of the  $\ell^1$ -norm recovers the same solution as minimization of the  $\ell^0$ -“norm”), including the motivation behind this relationship and a brief tour of various equivalence theorems in this field.

#### 9.1. Motivation from Compressed Sensing

Suppose that we wish to collect information about (i.e., sample or take measurements of) a continuous signal  $f(t)$  and then send or store this information in an efficient manner. For example,  $f(t)$  could be a sound wave or an image. Also suppose that a good approximation of the original signal must later be recovered. According to the Nyquist/Shannon sampling theorem, we must sample  $f(t)$  at a rate of at least twice its maximum frequency in order to be able to reconstruct  $f(t)$  exactly [82]. But in some applications, doing so may be expensive or even impossible.

In the circumstances that we are able to take many measurements of  $f(t)$  to obtain its discrete analog  $\mathbf{f} \in \mathbb{R}^N$ , one efficient method of compressing it is the following procedure: Let the columns of  $\Psi := [\psi_1, \dots, \psi_N]$  form an orthonormal basis for  $\mathbb{R}^N$ , and suppose that  $\mathbf{f}$  has a sparse representation in this basis, i.e., that we can write  $\mathbf{f} = \sum_{j=1}^N \alpha_j \psi_j$ , where  $\alpha_j := \langle \mathbf{f}, \psi_j \rangle$ ,  $1 \leq j \leq N$ , and  $\boldsymbol{\alpha} := [\alpha_1, \dots, \alpha_N]^\top$  is sparse. Setting all but the  $k$  largest (in absolute value) entries of  $\boldsymbol{\alpha}$  to 0 in order to obtain  $\boldsymbol{\alpha}_k$ , it can be shown that  $\Psi \boldsymbol{\alpha}_k$  gives the best  $k$ -term least squares approximation of  $\mathbf{f}$  in this basis. Clearly, the sparser  $\boldsymbol{\alpha}$  is, the better approximation we will obtain of  $\mathbf{f}$ , and in the case that  $\|\boldsymbol{\alpha}\|_0 \leq k$ , we recover the exact solution. This is the basic idea behind JPEG compression, which uses the discrete cosine transform as the sparsifying basis  $\Psi$  [72].

## 9.1. MOTIVATION FROM COMPRESSED SENSING

---

The problem with this procedure is that it is inefficient to spend time and money collecting all  $N$  samples if we are only going to throw most (all but  $k$ ) of them away when the signal is compressed. This is the motivation behind *compressed sensing*, originally proposed by Candès and Tao [14] and Donoho [28] (see also Candès and Tao's work [15] and the paper by Candès et al. [13]). Let  $\Phi \in \mathbb{R}^{m \times N}$  be a *sensing* or *measurement* matrix with  $m < N$  and consider the underdetermined system

$$\mathbf{y}_0 := \Phi \mathbf{f} = \Phi \Psi \boldsymbol{\alpha} = X \boldsymbol{\alpha}$$

for sparse  $\boldsymbol{\alpha}$ , where we have set  $X := \Phi \Psi$ . Ideally, we would recover  $\mathbf{f}$  by solving the optimization problem

$$(9.1) \quad \boldsymbol{\alpha}_0 := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \|\boldsymbol{\alpha}\|_0 \text{ subject to } X \boldsymbol{\alpha} = \mathbf{y}_0,$$

and set  $\hat{\mathbf{f}} := \Psi \boldsymbol{\alpha}_0$  with  $\hat{\mathbf{f}} \approx \mathbf{f}$ . However, solving Eq. (9.1) is NP-hard. Fortunately, when  $X$  satisfies certain conditions and when  $\boldsymbol{\alpha}_0$  is sufficiently sparse, the solution to Eq. (9.1) can be found by solving the  $\ell^1$ -minimization problem

$$(9.2) \quad \boldsymbol{\alpha}_1 := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \|\boldsymbol{\alpha}\|_1 \text{ subject to } X \boldsymbol{\alpha} = \mathbf{y}_0.$$

This was a riveting finding, as the optimization problem in Eq. (9.2) is convex and can be solved efficiently. It has been shown that this procedure (under certain conditions, e.g., when the columns of  $\Phi$  are uniformly random on the sphere  $S^{m-1}$ ) produces an approximation of  $\mathbf{f}$  that is as good as that of its best  $k$ -term approximation [28]. Further, theoretical and experimental results demonstrate that in many situations, the number of measurements  $m$  needed to recover  $\mathbf{f}$  is significantly less than  $N$  and can be much lower than the number required by the Nyquist/Shannon theorem. For example, when the measurement matrix  $\Phi \in \mathbb{R}^{m \times N}$  contains i.i.d. Gaussian entries, then exact recovery of  $\boldsymbol{\alpha}_0$  via  $\ell^1$ -minimization can be achieved (with high probability) in only  $m = O(k \log(N/k))$  measurements, where  $\|\boldsymbol{\alpha}_0\|_0 = k$  [14].

Even more astoundingly, similar results hold in the presence of noise, in which case the noiseless vector  $\mathbf{y}_0$  is replaced with  $\mathbf{y} = \mathbf{y}_0 + \mathbf{z}$ , where  $\mathbf{z} \in \mathbb{R}^m$  is a vector of errors satisfying  $\|\mathbf{z}\|_2 \leq \epsilon$ . Under certain conditions (involving how similar  $X$  is to an orthonormal basis when applied to

## 9.2. REVIEW OF EQUIVALENCE GUARANTEES

---

sufficiently-sparse vectors—see Definition 9.2.1), the  $\ell^1$ -minimization problem

$$(9.3) \quad \boldsymbol{\alpha}_{1,\epsilon} := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \|\boldsymbol{\alpha}\|_1 \text{ subject to } \|X\boldsymbol{\alpha} - \mathbf{y}\|_2 \leq \epsilon$$

is guaranteed to recover a coefficient vector approximating the ground truth sparse vector  $\boldsymbol{\alpha}_0$  (the solution to Eq. (9.1)) with  $\|\boldsymbol{\alpha}_{1,\epsilon} - \boldsymbol{\alpha}_0\|_2 \leq C_k \epsilon$  [13]. The constant  $C_k$  depends on properties of the matrix  $X$ .

A popular application of compressed sensing is *magnetic resonance imaging* (MRI), in which the measurement matrix  $\Phi$  consists of  $m$  randomly-selected rows of the discrete Fourier transform in  $\mathbb{R}^{N \times N}$  [64]. Other applications abound in the areas of data acquisition and compression, including sensor networks [105], seismology [51], and single pixel cameras [34].

The conditions under which the solutions to Eq. (9.1) and Eq. (9.2) are equal, i.e.,  $\ell^1/\ell^0$ -equivalence holds, involve the number of nonzero coordinates of the sparsest solution  $\boldsymbol{\alpha}_0$  and attributes of the matrix  $X$ . In the next section, we review several recovery guarantees, i.e., conditions guaranteeing exact or approximate recovery of the sparsest solution via  $\ell^1$ -minimization. To make the problem more general, we no longer explicitly assume the use of a sparsifying transform matrix  $\Psi$  and consider the general system  $X\boldsymbol{\alpha} = \mathbf{y}_0$ , for  $X \in \mathbb{R}^{m \times N}$  with  $m < N$ .

### 9.2. Review of Equivalence Guarantees

Existing results concerning  $\ell^1/\ell^0$ -equivalence consider the *incoherence* (or *spread*) of the vectors in the dictionary. Essentially, these results cannot be applied when the vectors are too correlated. A prototypical example is that if the data set contains two copies of the same vector (i.e., a pair of maximally-correlated vectors), then the minimum  $\ell^1$ -norm solution may contain a nonzero coefficient at either one of the copies or at a combination of the two. Contrast this with the sparsest solution, which would never contain nonzero coefficients at both copies.

There are various ways of measuring the incoherence in a dictionary, each leading its own theory relating the solutions of Eq. (9.1) and Eq. (9.2) (or its noise version Eq. (9.3)). We review a few approaches here.

## 9.2. REVIEW OF EQUIVALENCE GUARANTEES

---

### 9.2.1. Restricted Isometry.

DEFINITION 9.2.1. *For any integer  $k \in \{1, \dots, N\}$ , define the isometry constant of a matrix  $X \in \mathbb{R}^{m \times N}$  as the smallest number  $\delta_k$  such that*

$$(1 - \delta_k) \|\boldsymbol{\alpha}\|_2^2 \leq \|X\boldsymbol{\alpha}\|_2^2 \leq (1 + \delta_k) \|\boldsymbol{\alpha}\|_2^2$$

*for all vectors  $\boldsymbol{\alpha} \in \mathbb{R}^N$  satisfying  $\|\boldsymbol{\alpha}\|_0 \leq k$ .*

Put simply, the isometry constant  $\delta_k$  is a quantification of how close any set of  $k$  columns of  $X$  is to being an orthonormal basis. The following theorem uses the definition of isometry constants to relate the sparsest solution and that found by  $\ell^1$ -minimization:

THEOREM 9.2.1 (Cai, Wang, and Xu [9], Theorem 3.2). *For  $k > 1$ , suppose that  $\mathbf{y} = \mathbf{y}_0 + \mathbf{z}$ , where  $\|\mathbf{z}\|_2 \leq \epsilon$ ,  $\mathbf{y}_0 = X\boldsymbol{\alpha}$ , and  $\|\boldsymbol{\alpha}\|_0 \leq k$ . If the isometry constant  $\delta_k$  of  $X$  satisfies  $\delta_k < 0.307$ , then*

$$\|\boldsymbol{\alpha}_{1,\epsilon} - \boldsymbol{\alpha}\|_2 \leq \frac{\epsilon}{0.307 - \delta_k},$$

*where  $\boldsymbol{\alpha}_{1,\epsilon}$  is the solution to Eq. (9.3).*

Note that since this holds for all  $\boldsymbol{\alpha}$  satisfying  $\mathbf{y}_0 = X\boldsymbol{\alpha}$  with  $\|\boldsymbol{\alpha}\|_0 \leq k$ , it must be that  $\boldsymbol{\alpha} = \boldsymbol{\alpha}_0$  is the sparsest solution. In the case that  $\epsilon = 0$ , we have that  $\boldsymbol{\alpha}_{1,\epsilon} = \boldsymbol{\alpha} = \boldsymbol{\alpha}_0$ , and so  $\ell^1/\ell^0$ -equivalence holds! Further, this condition on  $\delta_k$  is essentially tight; there exists a case for which  $\ell^1/\ell^0$ -equivalence does *not* hold and  $\delta_k = (k-1)/(2k-1)$  [9]; considering that  $(k-1)/(2k-1) = 1/3$  for  $k = 2$  and  $(k-1)/(2k-1) < 0.5$  for all  $k \in \mathbb{N}$ , this does not leave much room for improvement.

Random matrices, e.g., Gaussian and Bernoulli ensembles, have particularly small restricted isometry constants. For example, when the entries of  $X$  are drawn from  $\mathcal{N}(0, 1/m)$ ,  $\ell^1/\ell^0$ -equivalence holds with overwhelmingly high probability provided that, for  $\|\boldsymbol{\alpha}_0\|_0 \leq k$ , the ratio  $k/N$  is sufficiently small. [14].

In the case that the coefficient vector  $\boldsymbol{\alpha}$  is not sparse, the definition of restricted isometry can also be used to show that the  $\ell^1$ -minimized solution is close to  $\boldsymbol{\alpha}$ , provided that  $\boldsymbol{\alpha}$  has a good  $k$ -term approximation:

## 9.2. REVIEW OF EQUIVALENCE GUARANTEES

---

THEOREM 9.2.2 (Cai, Wang, and Xu [9], Theorem 3.3). *Given any noisy approximation  $\mathbf{y} = \mathbf{y}_0 + \mathbf{z}$  with  $\mathbf{y}_0 = X\boldsymbol{\alpha}$ , where  $\|\mathbf{z}\|_2 \leq \epsilon$  and the isometry constant  $\delta_k$  of  $X$  satisfies  $\delta_k < 0.307$ , then*

$$\|\boldsymbol{\alpha}_{1,\epsilon} - \boldsymbol{\alpha}\|_2 \leq \frac{\epsilon}{0.307 - \delta_k} + \frac{1}{0.307 - \delta_k} \frac{\|\boldsymbol{\alpha} - \boldsymbol{\alpha}_k\|_1}{\sqrt{k}},$$

where  $\boldsymbol{\alpha}_{1,\epsilon}$  is the solution to Eq. (9.3) and  $\boldsymbol{\alpha}_k$  is the result of setting all but the  $k$  largest entries (in magnitude) of  $\boldsymbol{\alpha}$  to 0.

When  $\boldsymbol{\alpha}$  satisfies  $\|\boldsymbol{\alpha}\|_0 \leq k$ , this reduces to Theorem 9.2.1.

Despite the wide applicability of this approach, computing the isometry constant to verify the conditions of the above theorems is often difficult. To calculate  $\delta_k$  for a given  $X$ , we must examine all  $k$ -element subsets of the columns of  $X$ , which is expensive, prohibitively so when  $N$  is large. Unfortunately, this is also the case with the next notion of incoherence we examine.

### 9.2.2. Spark.

DEFINITION 9.2.2. *The spark of a matrix  $X$ , denoted  $\text{Spark}(X)$ , is the smallest number of linearly dependent columns of  $X$ .*

Though it initially appears quite similar to matrix rank,  $\text{Spark}(X)$  is combinatorial to compute, requiring the examination of all increasingly large subsets of columns of  $X$  until a linearly-dependent subset is found. However, if it can be determined, the following theorem can be applied:

THEOREM 9.2.3 (Donoho and Elad [30], Corollary 4). *If a representation  $X\boldsymbol{\alpha} = \mathbf{y}_0$  satisfies  $\|\boldsymbol{\alpha}\|_0 < \text{Spark}(X)/2$ , then  $\boldsymbol{\alpha}$  is necessarily the sparsest such representation, i.e.,  $\boldsymbol{\alpha} = \boldsymbol{\alpha}_0$ , the solution to Eq. (9.1).*

The relevant consequence of Theorem 9.2.3 is that if  $\ell^1$ -minimization finds a coefficient vector  $\boldsymbol{\alpha}_1$  (the solution to Eq. (9.2)) such that  $\|\boldsymbol{\alpha}_1\|_0 < \text{Spark}(X)/2$ , then  $\ell^1/\ell^0$ -equivalence holds. This bound is sharp in the sense that no quantity smaller than  $\text{Spark}(X)/2$  can replace it in Theorem 9.2.3 [30].

To illustrate this theorem, consider that a random  $m \times N$  matrix with full rank has spark equal to  $m + 1$  with high probability, and so we only need  $\|\boldsymbol{\alpha}_1\|_0 \leq (m + 1)/2$  to confirm that it is the sparsest solution.

## 9.2. REVIEW OF EQUIVALENCE GUARANTEES

---

Despite the difficulty in computing  $\text{Spark}(X)$ , one can see that if there were a tractable lower bound on this quantity, Theorem 9.2.3 would become more useful. This brings us to the notion of *mutual coherence*.

### 9.2.3. Mutual Coherence.

**DEFINITION 9.2.3.** *Given a matrix  $X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{m \times N}$  with normalized columns (so that  $\|\mathbf{x}_j\|_2 = 1$  for  $1 \leq j \leq N$ ), the mutual coherence of  $X$ , denoted  $\mu(X)$ , is given by*

$$\mu(X) = \max_{1 \leq i \neq j \leq N} |\langle \mathbf{x}_i, \mathbf{x}_j \rangle|.$$

Note that mutual coherence is straightforward to compute for a given dictionary. The relationship  $\text{Spark}(X) \geq 1 + 1/\mu(X)$  [8, 30] produces the following theorem:

**THEOREM 9.2.4** (Donoho and Elad [30], Theorem 12; Gribonval and Nielsen [46], Theorem 1). *Let  $X \in \mathbb{R}^{m \times N}$ ,  $m < N$ , have normalized columns and mutual coherence  $\mu(X)$ . If  $\boldsymbol{\alpha}$  satisfies  $X\boldsymbol{\alpha} = \mathbf{y}_0$  with*

$$(9.4) \quad \|\boldsymbol{\alpha}\|_0 < \frac{1}{2} \left( 1 + \frac{1}{\mu(X)} \right),$$

*then  $\boldsymbol{\alpha}$  is the unique solution to the  $\ell^1$ -minimization problem in Eq. (9.2).*

Again, this means that if  $\ell^1$ -minimization finds a solution with less than  $(1/2)(1 + \mu(X)^{-1})$  nonzeros, then it is necessarily the sparsest solution and so  $\ell^1/\ell^0$ -equivalence holds.

Let us differentiate the noise tolerance  $\zeta$  from the approximation error bound  $\epsilon$  in Eq. (9.3), so that we can bound  $\|\mathbf{z}\|_2 \leq \zeta$  and  $\|X\boldsymbol{\alpha} - \mathbf{y}\|_2 \leq \epsilon$  separately.<sup>1</sup> (Note that  $\zeta = \epsilon$  in Theorems 9.2.1 and 9.2.2.) Then the following theorem by Donoho et al. gives conditions for  $\ell^1/\ell^0$ -equivalence in the noisy setting:

**THEOREM 9.2.5** (Donoho, Elad, and Temlyakov [31], Theorem 3.1). *Let  $X \in \mathbb{R}^{m \times N}$ ,  $m < N$ , have normalized columns and mutual coherence  $\mu(X)$ . Suppose there exists an ideal noiseless signal*

---

<sup>1</sup>Given that  $\mathbf{y} = \mathbf{y}_0 + \mathbf{z}$  with  $\|\mathbf{z}\|_2 \leq \zeta$ , we have that  $\|X\boldsymbol{\alpha} - \mathbf{y}\|_2 \leq \|X\boldsymbol{\alpha} - \mathbf{y}_0\|_2 + \zeta$  for any coefficient vector  $\boldsymbol{\alpha}$ . This upper bound may be larger than  $\zeta$  when  $\boldsymbol{\alpha} \neq \boldsymbol{\alpha}_0$ , and so we may set  $\epsilon \geq \zeta$ , as in the experiments in Chapter 11. However, since the bound  $\|X\boldsymbol{\alpha} - \mathbf{y}\|_2 \leq \|X\boldsymbol{\alpha} - \mathbf{y}_0\| + \|\mathbf{z}\|_2$  could be very loose, we could also set  $\epsilon < \zeta$ .

## 9.2. REVIEW OF EQUIVALENCE GUARANTEES

---

$\mathbf{y}_0$  such that  $\mathbf{y}_0 = X\boldsymbol{\alpha}$  and

$$(9.5) \quad \|\boldsymbol{\alpha}\|_0 = k \leq \frac{1}{4} \left( 1 + \frac{1}{\mu(X)} \right).$$

Then  $\boldsymbol{\alpha} = \boldsymbol{\alpha}_0$  is the unique sparsest representation of  $\mathbf{y}_0$  over  $X$ . Further, suppose that we only observe  $\mathbf{y} = \mathbf{y}_0 + \mathbf{z}$  with  $\|\mathbf{z}\|_2 \leq \zeta$ . Then we have

$$(9.6) \quad \|\boldsymbol{\alpha}_{1,\epsilon} - \boldsymbol{\alpha}_0\|_2^2 \leq \frac{(\epsilon + \zeta)^2}{1 - \mu(X)(4k - 1)},$$

where  $\boldsymbol{\alpha}_{1,\epsilon}$  is the solution to Eq. (9.3).

That is, if the ideal sparse vector  $\boldsymbol{\alpha}_0$  is sparse enough and the mutual coherence of  $X$  is small enough,  $\ell^1$ -minimization will give us a solution close to  $\boldsymbol{\alpha}_0$ , with “how close” depending on  $k$ ,  $\mu(X)$ , the noise tolerance  $\zeta$ , and the approximation error bound  $\epsilon$ .

We can also say something regarding the support of  $\boldsymbol{\alpha}_{1,\epsilon}$  in the noisy setting:

**THEOREM 9.2.6** (Donoho, Elad, Temlyakov [31], Theorem 6.1). *Suppose that  $\mathbf{y} = \mathbf{y}_0 + \mathbf{z}$ , where  $\mathbf{y}_0 = X\boldsymbol{\alpha}_0$ ,  $\|\boldsymbol{\alpha}\|_0 \leq k$  and  $\|\mathbf{z}\|_2 \leq \zeta$ . Suppose that  $\beta := \mu(X)k < \frac{1}{2}$  (so  $k < \frac{1}{2\mu(X)}$ ). Set*

$$(9.7) \quad \gamma := \frac{\sqrt{1 - \beta}}{1 - 2\beta}.$$

*Then given  $\boldsymbol{\alpha}_{1,\epsilon}$  the solution to Eq. (9.3) with exaggerated error tolerance  $\epsilon := C\zeta$  where  $C = C(\mu(X), k) := \gamma\sqrt{k}$ , we have that  $\text{supp}(\boldsymbol{\alpha}_{1,\epsilon}) \subset \text{supp}(\boldsymbol{\alpha}_0)$ .*

This says that when the mutual coherence is very small relative to the sparsity level, the solution  $\boldsymbol{\alpha}_{1,\epsilon}$  to Eq. (9.3) has the same support as the sparsest solution  $\boldsymbol{\alpha}_0$ . (Observe that  $\boldsymbol{\alpha}_0$  is indeed the sparsest solution by Theorem 9.2.4, since  $\|\boldsymbol{\alpha}_0\|_0 < (1/2)\mu(X)^{-1} < (1/2)(1 + \mu(X)^{-1})$ .) Since  $\epsilon = \gamma\sqrt{k}\zeta$  and  $\gamma \geq 1$ ,  $\epsilon \geq \zeta$  is required in Theorem 9.2.6.

Despite the ease of computing mutual coherence, these theorems produce what are generally considered to be fairly loose bounds on the sparsity level  $k$ , given experimental results and cases for which restricted isometry constants are known [49]. As an example, consider that  $\text{Spark}(X)$  can be as large as  $m + 1$ , whereas  $1 + \mu(X)^{-1}$  is bounded above by  $1 + \sqrt{m}$ , assuming that the columns of  $X$  are normalized. Thus we expect that Theorem 9.2.3 will generally allow us to take many less measurements (and still prove  $\ell^1/\ell^0$ -equivalence) than Theorem 9.2.4.

## 9.2. REVIEW OF EQUIVALENCE GUARANTEES

---

There are other recovery guarantees involving incoherence that hold with high probability. When applied to random matrices, these guarantees are generally stronger than those discussed above (in terms of requiring less measurements and/or less sparsity of the solution vector). For example, Candès and Plan [11] provide conditions that guarantee recovery (with high probability) of sparse and approximately sparse solutions in the case that the rows of the dictionary are sampled independently from certain probability distributions. These conditions are in terms of incoherence, defined as an upper bound on the squared norms of the rows of  $X$  (either deterministically or stochastically), and require an *isotropy* property [11]. In the case that the probability distribution has mean 0, this property states that the covariance matrix of the probability distribution is equal to the identity matrix.

In another paper [10], Candès and Plan guarantee probabilistic recovery in terms of a condition on mutual coherence (as defined in Definition 9.2.3) that is satisfied with high probability on certain random matrices. These recovery guarantees allow for the sparsity level  $k$  in the case of these random matrices to be notably larger than in Eq. (9.4) in Theorem 9.2.4. We also mention the results by Tropp [92] concerning recovery in terms of mutual coherence and the extreme singular values of randomly-chosen subsets of dictionary columns.

**9.2.4. Neighborliness.** Thus far, the equivalence guarantees we have looked at have given only *sufficient* conditions for  $\ell^1/\ell^0$ -recovery. The following notion of incoherence provides a condition that is also *necessary*. Define  $P$ , the polytope associated with the matrix  $X \in \mathbb{R}^{m \times N}$ , as  $P := XC$ , where  $C$  is the  $N$ -dimensional polytope determined by  $C := \{\mathbf{x} \in \mathbb{R}^N : \|\mathbf{x}\|_1 \leq 1\}$ .

**DEFINITION 9.2.4.** *The polytope  $P$  associated with the matrix  $X$  is said to be  $k$ -neighborly if every set of  $k + 1$  vertices (that does not include an antipodal pair) forms a face of  $P$ .*

**THEOREM 9.2.7** (Donoho [27], Theorem 1 (see also Donoho and Tanner [32])). *The polytope  $P$  associated with the matrix  $X$  has  $2N$  vertices and is  $k$ -neighborly if and only if whenever the sparsest solution  $\boldsymbol{\alpha}_0$  satisfying  $X\boldsymbol{\alpha}_0 = \mathbf{y}$  has no more than  $k$  nonzeros,  $\ell^1$ -minimization in the form of Eq. (9.2) recovers  $\boldsymbol{\alpha}_0$ .*

That is, the associated polytope having  $2N$  vertices and being  $k$ -neighborly is a necessary and sufficient condition for  $\ell^1/\ell^0$ -equivalence in the case that  $\|\boldsymbol{\alpha}_0\|_0 \leq k$ . In general, polytopes of the

## 9.2. REVIEW OF EQUIVALENCE GUARANTEES

---

form  $P$  are no more than  $k \leq \lceil (m+1)/3 \rceil$  neighborly [27]. Combining this fact with Theorem 9.2.7 produces the following corollary:

**COROLLARY 9.2.1** (Donoho [27], Corollary 1.3 (see also Donoho and Tanner [32])). *For  $N-2 \geq m > 2$ , if  $\ell^1$ -minimization in the form of Eq. (9.2) recovers all sparse vectors satisfying  $\|\alpha\|_0 \leq k$ , then  $k \leq \lceil (m+1)/3 \rceil$ .*

Unfortunately, determining whether  $P$  is  $k$ -neighborly is not tractable for large dimension  $m$ . In the SRC paper [104], Wright et al. used a combinatorial (but relatively efficient) algorithm [83] to compute a tighter upper bound on polytope neighborliness and used this bound to estimate for which dimension (i.e., number of measurements)  $m$   $\ell^1/\ell^0$ -equivalence can be achieved. However, this is a very imprecise approach, as an upper bound on neighborliness can only provide us with a negative result of the following form: If  $P$  is no more than  $k'$ -neighborly and  $\|\alpha_1\|_0 := k > k'$ , then  $\alpha_1 \neq \alpha_0$ , i.e.,  $\ell^1/\ell^0$ -equivalence does not hold. However, in the case that  $k \leq k'$ , nothing more can be said.

## CHAPTER 10

### Mutual Coherence Equivalence in the Context of Classification

In classification problems, data from the same class may be highly-correlated. As we will show, the  $\ell^1/\ell^0$ -recovery guarantees that can be applied in practice often do not hold. Thus the “theory” behind sparse representation-based methods for classification is missing a significant piece. Though such methods work experimentally well, is their success really due to sparsity? Is sparsity really being achieved? We aim to provide insight into these questions via a series of projects in the next few chapters.

Though in the last chapter we used  $\mathbf{y}_0$  to refer to a clean measurement vector and  $\mathbf{y} := \mathbf{y}_0 + \mathbf{z}$  to refer to its noisy version, in this chapter and in Chapter 12,  $\mathbf{y}$  may represent *either* a clean or noisy measurement vector, or more commonly, an arbitrary test sample (as was done in Part 1). We do this because, in the context of representation-based classification, there are reasons other than noise in the test sample for allowing the equality  $\mathbf{y} = X_{\text{tr}}\boldsymbol{\alpha}$  to hold only approximately: the training data could also be corrupted, or we may want to relax the assumption that class manifolds are linear subspaces (perhaps this is only approximately, or locally, the case). To keep the situation general and to avoid confusion, we will only differentiate between  $\mathbf{y}$  and  $\mathbf{y}_0$  when we explicitly consider  $\mathbf{z}$  as noise, as in Donoho et al.’s Theorems 9.2.5 and 9.2.6 that consider mutual coherence in the case of noise.

That said, it is important to note that in the case that the equality  $\mathbf{y} = X_{\text{tr}}\boldsymbol{\alpha}$  holds only approximately due to class manifold curvature, theorems such as Theorem 9.2.5 can still (in theory) be applied to show that the solution found by  $\ell^1$ -minimization approximates the sparsest solution. Here, the vector  $\mathbf{z}$  can refer to the error between a *clean* test sample  $\mathbf{y}$  and the projection of  $\mathbf{y}$  onto the span of (a portion of) its same-class training samples. In this case, the notation  $\hat{\mathbf{y}}_l$  for the projection of a class  $l$  test sample  $\mathbf{y}$  onto the span of local class  $l$  training samples, with  $\mathbf{y} = \hat{\mathbf{y}}_l + \mathbf{z}$ , is perhaps more appropriate than using  $\mathbf{y} = \mathbf{y}_0 + \mathbf{z}$ . This is another reason to generalize our notation.

## 10.1. THE DILEMMA

---

### 10.1. The Dilemma

In this chapter, we investigate whether or not  $\ell^1/\ell^0$ -recovery guarantees involving mutual coherence are compatible with the assumptions made in SRC. As we saw in the last chapter, this type of recovery guarantee is easily computable, unlike those involving restricted isometry constants, spark, or polytope neighborliness.

Ideally, we could use the mutual coherence conditions given in Theorem 9.2.4 or 9.2.5 to show that  $\ell^1$ -minimization in SRC produces a solution equal or close to the sparsest solution. However, in classification problems, samples in the same class can be extremely correlated, resulting in  $\mu(X_{\text{tr}})$  being close to 1 (we assume that the training samples have  $\ell^2$ -norm equal to 1, throughout). For example, recall the sinusoidal wave synthetic database used in Chapter 5: With training class size  $N_0 = 25$  and noise level  $\eta = 0.001$ , the average mutual coherence of  $X_{\text{tr}}$  over 100 trials is 1.0000. As further evidence of the large mutual coherence inherent in classification problems, Table 10.1 displays the average values of  $\mu(X_{\text{tr}})$  (over 10 trials) of the face databases used in Chapter 5. The training sets reflect 50/50 training/test splits.

Database	$m_{\text{PCA}} = 30$	$m_{\text{PCA}} = 56$	$m_{\text{PCA}} = 120$
AR-1	0.9991	0.9987	0.9985
AR-2	0.9993	0.9988	0.9984
Ext Yale B	0.9951	0.9954	0.9941
ORL	0.9971	0.9970	0.9966

TABLE 10.1. Average mutual coherence (over 10 trials) computed from the samples in  $X_{\text{tr}}$  (50/50 training/test splits) of the face databases from Chapter 5.

When  $\mu(X_{\text{tr}}) \approx 1$ , the mutual coherence bound in Theorem 9.2.4 becomes

$$\|\boldsymbol{\alpha}\|_0 < \frac{1}{2} \left( 1 + \frac{1}{\mu(X_{\text{tr}})} \right) \approx 1.$$

Since  $\|\boldsymbol{\alpha}\|_0$  denotes the number of nonzero coefficients in the representation of  $\mathbf{y}$  over  $X_{\text{tr}}$ , it will never satisfy  $\|\boldsymbol{\alpha}\|_0 < 1$ . Thus we cannot use Theorem 9.2.4 to prove  $\ell^1/\ell^0$ -equivalence in SRC, for example, on the databases used in Chapter 5.

### 10.3. PRELIMINARY RESULTS

---

#### 10.2. Main Goal

In the following sections, we identify cases in which the condition given in Eq. (9.4) from Theorem 9.2.4 *provably* does not hold, and thus we cannot use Theorem 9.2.4 to prove  $\ell^1/\ell^0$ -equivalence. We also discuss analogous results in the noisy case, i.e., Eq. (9.5) in Theorem 9.2.5. In particular, we are concerned with the applicability of these theorems for class-structured data.

#### 10.3. Preliminary Results

We will use the following lemma which gives a lower-bound on mutual coherence in the under-determined setting:

LEMMA 10.3.1 (Rosenfeld [77], Theorem 3). *For  $X \in \mathbb{R}^{m \times N}$  with normalized columns and  $m < N$ , we have that*

$$\mu(X) \geq \sqrt{\frac{N-m}{m(N-1)}}.$$

It is straightforward to show that Lemma 10.3.1 implies that

$$\mu(X) \geq \frac{1}{m},$$

since  $\sqrt{\frac{N-m}{m(N-1)}}$  monotonically increases in  $N \in \mathbb{N}$  for  $N > m$ , with a minimum value of  $1/m$  attained at  $N = m + 1$ . Thus to have even a *chance* of Theorem 9.2.4 or 9.2.5 holding, we must have

$$\|\boldsymbol{\alpha}\|_0 < \frac{1}{c} \left( 1 + \frac{1}{\mu(X)} \right) \leq \frac{1}{c} (1 + m),$$

where  $c = 2$  in the noiseless case and  $c = 4$  in the noisy setting.

We next consider the smallest possible value of the number of nonzeros  $\|\boldsymbol{\alpha}\|_0$  in any representation  $X\boldsymbol{\alpha} = \mathbf{y}$ . Let us assume that each test sample is not a scalar multiple of any training sample. It follows that  $\|\boldsymbol{\alpha}\|_0$  is at least 2. Thus in order for Theorem 9.2.4 or 9.2.5 to hold, we must have

$$2 \leq \|\boldsymbol{\alpha}\|_0 < \frac{1}{c} \left( 1 + \frac{1}{\mu(X_{\text{tr}})} \right) \Rightarrow \mu(X_{\text{tr}}) < \frac{1}{2c-1}$$

## 10.4. MAIN RESULT

---

$$\Rightarrow \mu(X_{\text{tr}}) < \begin{cases} 1/3, & \text{noiseless case} \\ 1/7, & \text{noisy setting.} \end{cases}$$

Note that these upper bounds for  $\mu(X_{\text{tr}})$  are very small compared to the values of  $\mu(X_{\text{tr}})$  in Table 10.1!

These findings produce the following small-scale result:

**PROPOSITION 10.3.1.** *Suppose that  $X_{\text{tr}}\alpha = \mathbf{y}$ . If  $m \leq 3$  and  $\mathbf{y}$  is not a scalar multiple of any training sample, then the inequality in Eq. (9.4) with  $X = X_{\text{tr}}$  does not hold. That is, we cannot use Theorem 9.2.4 to prove  $\ell^1/\ell^0$ -equivalence in SRC.*

**Proof.** By Lemma 10.3.1, we must have that  $\mu(X_{\text{tr}}) \geq \frac{1}{m} \geq \frac{1}{3}$ .  $\square$

An analogous statement holds in the noisy setting (Theorem 9.2.5) for  $m \leq 7$ .

### 10.4. Main Result

**PROPOSITION 10.4.1 (Main Result).** *Suppose that the sparsest representation of  $\mathbf{y} \in \mathbb{R}^m$  over the dictionary  $X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{m \times N}$  is given by  $\mathbf{y} = \alpha_1 \mathbf{x}_{j_1} + \dots + \alpha_{j_k} \mathbf{x}_{j_k}$  for  $\{j_1, \dots, j_k\} \subset \{1, \dots, N\}$ . Set  $\tilde{N}$  to be the number of columns of  $X$  contained in*

$$\tilde{\mathcal{X}} := \text{span}\{\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_k}\}.$$

*Of course,  $\tilde{N} \geq k$ . If  $\tilde{N} > k$ , then the inequality in Eq. (9.4) does not hold. That is, we cannot use Theorem 9.2.4 to prove  $\ell^1/\ell^0$ -equivalence.*

**Proof:** Suppose that  $\tilde{N} > k$ . Then there are more than  $k$  dictionary elements in the subspace  $\tilde{\mathcal{X}}$ . Since the vectors  $\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_k}$  are linearly independent (because otherwise,  $\mathbf{y}$  could be expressed more sparsely), the dimension of  $\tilde{\mathcal{X}}$  is exactly  $k$ .

Define  $\tilde{X} \in \mathbb{R}^{m \times \tilde{N}}$  to be the matrix of the  $\tilde{N}$  dictionary elements contained in  $\tilde{\mathcal{X}}$ . Let the singular value decomposition of  $\tilde{X}$  be given by  $\tilde{X} = U\Sigma V^T$ , and set  $U_k$  to contain the first  $k$  columns of  $U$ ,  $V_k$  to contain the first  $k$  columns of  $V$ , and  $\Sigma_k$  to contain the first  $k$  columns and rows of  $\Sigma$ . Because  $\tilde{X}$  has rank  $k$ , we can alternatively write

$$\tilde{X} = U_k \Sigma_k V_k^T.$$

## 10.4. MAIN RESULT

---

The  $k \times \tilde{N}$  matrix  $U_k^\top \tilde{X}$  has the same mutual coherence as  $\tilde{X}$ , since they have the same Gram matrices:

$$\begin{aligned}
(U_k^\top \tilde{X})^\top (U_k^\top \tilde{X}) &= \tilde{X}^\top U_k U_k^\top \tilde{X} \\
&= (U_k \Sigma_k V_k^\top)^\top U_k U_k^\top (U_k \Sigma_k V_k^\top) \\
&= V_k \Sigma_k^\top U_k^\top U_k U_k^\top U_k \Sigma_k V_k^\top \\
&= V_k \Sigma_k^\top \Sigma_k V_k^\top \\
&= (U_k \Sigma_k V_k^\top)^\top (U_k \Sigma_k V_k^\top) \\
&= \tilde{X}^\top \tilde{X}.
\end{aligned}$$

By Lemma 10.3.1, we have that

$$\mu(X) \geq \mu(\tilde{X}) = \mu(U_k^\top \tilde{X}) \geq \sqrt{\frac{\tilde{N} - k}{k(\tilde{N} - 1)}} \geq \sqrt{\frac{(k+1) - k}{k((k+1) - 1)}} = \frac{1}{k}.$$

Thus the bound on  $k$  in Theorem 9.2.4 requires that

$$(10.1) \quad k < \frac{1}{2} \left( 1 + \frac{1}{\mu(X)} \right) \leq \frac{1}{2}(1+k) \Rightarrow k < 1. \quad \square$$

We present several corollaries to Proposition 10.4.1. The first is a consequence applicable to any  $\ell^1$ -minimization problem, regardless of whether or not the dictionary elements have class structure:

**COROLLARY 10.4.1** (Consequence for general  $\ell^1$ -minimization). *If a measurement vector  $\mathbf{y} \in \mathbb{R}^m$  is not at all sparse over the dictionary  $X \in \mathbb{R}^{m \times N}$ , i.e., if every representation of  $\mathbf{y}$  requires no less than  $m$  dictionary elements, then the condition in Eq. (9.4) from Theorem 9.2.4 does not hold. This is because the dimension of  $\tilde{\mathcal{X}}$  (as defined in Proposition 10.4.1) is  $m$ , and so every dictionary element is contained in  $\tilde{\mathcal{X}}$ .*

Corollary 10.4.1 illustrates the importance of choosing a dictionary that awards a sparse representation of  $\mathbf{y}$  in any application of  $\ell^1$ -minimization.

The following corollary follows from the proof of Proposition 10.4.1:

## 10.4. MAIN RESULT

---

COROLLARY 10.4.2. *Let  $X \in \mathbb{R}^{m \times N}$  with  $m < N$ , and let  $k$  be any positive integer such that  $k < N$ . If any set of  $k$  linearly independent columns of  $X$  spans an additional, distinct column of  $X$ , then the bound*

$$k < \frac{1}{2} \left( 1 + \frac{1}{\mu(X)} \right)$$

*does not hold.*

Of course, this bound will not hold for any larger values of  $k$ , either. This means that if we can find an integer  $k$  satisfying the conditions of Corollary 10.4.2, then any attempt to prove  $\ell^1/\ell^0$ -equivalence using Theorem 9.2.4 will require  $X\alpha = \mathbf{y}$  with  $\|\alpha\|_0 < k$ . Recalling the previous chapter, the hypothesis of Corollary 10.4.2 states that  $\text{Spark}(X) \leq k + 1$ . Using that  $\text{Spark}(X) \geq 1 + 1/\mu(X)$  [8, 30], we can alternatively prove Corollary 10.4.2 by noting that

$$k \geq \frac{k+1}{2} \geq \frac{\text{Spark}(X)}{2} \geq \frac{1}{2} \left( 1 + \frac{1}{\mu(X)} \right).$$

Though computing  $\text{Spark}(X)$  is generally not tractable, Corollary 10.4.2 only requires computing an upper bound for  $\text{Spark}(X)$ . Further, if  $X = X_{\text{tr}}$  is the training set in a classification problem, we might expect that  $\text{Spark}(X_{\text{tr}})$  is lower than that for arbitrary dictionaries due to the dictionary's class structure. Thus we may only need to try a few small values for  $k$  in order to apply Corollary 10.4.2.

The following corollary is an explicit consequence for class-structured dictionaries:

COROLLARY 10.4.3 (Consequence for Class-Structured Dictionaries). *Suppose that  $\mathbf{y}$  is a test sample with  $\|\mathbf{y}\|_2 = 1$ , and define  $\mu := \mu(X_{\text{tr}})$ . If the addition of  $\mathbf{y}$  does not increase the mutual coherence, that is, if  $|\langle \mathbf{y}, \mathbf{x}_i \rangle| \leq \mu$  for all  $1 \leq i \leq N_{\text{tr}}$ , i.e.,  $\mu([\mathbf{y}, X_{\text{tr}}]) = \mu$ , then we cannot have both that (i)  $X_{\text{tr}}\alpha = \mathbf{y}$  and (ii)  $\|\alpha\|_0 < (1/2)(1 + (1/\mu(X_{\text{tr}})))$ .*

**Proof:** If we can write  $X_{\text{tr}}\alpha = \mathbf{y}$  for  $\|\alpha\|_0 =: k$ , then the  $k$  (linearly independent) training samples with nonzero coefficients in the representation span a  $k$ -dimensional subspace containing  $\mathbf{y}$ . Setting  $X = [\mathbf{y}, X_{\text{tr}}]$  in Corollary 10.4.2, we have that

$$k \not< \frac{1}{2} \left( 1 + \frac{1}{\mu(X)} \right) = \frac{1}{2} \left( 1 + \frac{1}{\mu} \right).$$

## 10.4. MAIN RESULT

---

On the other hand, if

$$k < \frac{1}{2} \left( 1 + \frac{1}{\mu(X)} \right) = \frac{1}{2} \left( 1 + \frac{1}{\mu} \right)$$

for some positive integer  $k < N_{\text{tr}}$ , then also by Corollary 10.4.2, it must be the case that  $\mathbf{y}$  is not contained in the subspace spanned by any  $k$  linearly independent distinct columns of  $X$ , i.e., columns of  $X_{\text{tr}}$ . Thus we cannot write  $X_{\text{tr}}\boldsymbol{\alpha} = \mathbf{y}$  for any  $\boldsymbol{\alpha}$  satisfying  $\|\boldsymbol{\alpha}\|_0 = k$ .  $\square$

It might initially seem that the hypothesis of Corollary 10.4.3 is unlikely to hold. However, if one assumes that the data is sampled randomly with test samples having the same distribution as the training samples in their ground truth classes, then the hypothesis that  $\mu([\mathbf{y}, X_{\text{tr}}]) = \mu(X_{\text{tr}})$  becomes much more probable. This actually becomes a critical problem for us in Chapter 12.

Our final corollary determines conditions under which the bound in Eq. (9.4) from Theorem 9.2.4 is theoretically incompatible with the explicit assumptions made in SRC. We review these assumptions briefly:

**ASSUMPTION 1** (Linear Subspaces). *The ground truth class manifolds of the given data set are linear subspaces.*

**ASSUMPTION 2** (Spanning Training Set). *The training matrix  $X_{\text{tr}}$  contains sufficient samples in each class to span the corresponding linear subspace.*

**COROLLARY 10.4.4** (Consequence for SRC). *Suppose that the SRC Assumptions 1 and 2 hold. Let  $\mathbf{y}$  have ground truth class  $l$ , and suppose that the number of class  $l$  training samples,  $N_l$ , is large, i.e.,  $N_l > d_l$ , for  $d_l$  the dimension of the linear subspace representing the class  $l$  manifold. Then there exists a test sample  $\mathbf{y}$  which requires the maximum number  $d_l$  of class  $l$  training samples to represent it. If this representation of  $\mathbf{y}$  is its sparsest representation over the dictionary  $X_{\text{tr}}$ , then the condition in Eq. (9.4) from Theorem 9.2.4 cannot hold. Thus we cannot use Theorem 9.2.4 to prove  $\ell^1/\ell^0$ -equivalence in SRC.*

Corollary 10.4.4 says that if we have a surplus of class  $l$  training samples (i.e., more than enough to span the class  $l$  subspace), then, provided that the “class representations” (representations of the test samples in terms of their ground truth classes) truly are the sparsest representations of

## 10.4. MAIN RESULT

---

the test samples over the training set (as argued by the SRC authors [104]), there will be some test samples for which Theorem 9.2.4 cannot hold. These test samples are exactly those requiring  $k = d_l$  class  $l$  training samples in their representations. In general, such test samples must exist; otherwise, the dimension of the class  $l$  subspace would be less than  $d_l$ . To reiterate, if everything we *want* to happen in SRC happens (large class sizes, sparse class representations), then we cannot consistently use Theorem 9.2.4 to prove  $\ell^1/\ell^0$ -equivalence.

On a more positive note, the assumptions in SRC make it possible to estimate whether or not the conditions of Proposition 10.4.1 hold. Though these conditions are difficult to check in general (if we knew the sparsest solution of  $\mathbf{y}$  over the dictionary, then we would not need to use  $\ell^1$ -minimization to find it!), the linear subspace assumption in SRC gives us a heuristic for doing so. We could potentially estimate the dimension of each class (using one of the methods discussed in the context of the LPCA-SRC algorithm in Chapter 4, for instance) and compare this with the number of training samples in that class. If the latter is larger than the former, then we expect that Theorem 9.2.4 cannot be applied for some test samples.

In typical applications, we must deal with noisy data. Thus we should consider the application of Theorem 9.2.5 instead of Theorem 9.2.4. But this is immediate: Since the mutual coherence condition is stricter in the case of noise, the consequences of Proposition 10.4.1 and the above corollaries hold whenever the conditions are assumed to hold on the clean version of the data. In particular, Theorem 9.2.5 requires the *existence* of a clean test sample  $\mathbf{y}_0$  (even if it is unknown to us) that satisfies  $X\boldsymbol{\alpha} = \mathbf{y}_0$  with  $\|\boldsymbol{\alpha}\|_0 \leq (1/4)(1 + (1/\mu(X)))$ . Under the hypothesis of Corollary 10.4.3 (setting  $\mathbf{y}_0 = \mathbf{y}$ ), such a  $\mathbf{y}_0$  cannot exist. An analogous result holds in the case that class manifolds contain curvature, using the framework discussed in the beginning of this chapter.

As a concluding remark, we stress that the mutual coherence conditions in Theorems 9.2.4 and 9.2.5 are sufficient, but not necessary, for  $\ell^1/\ell^0$ -equivalence. Thus it is possible for  $\ell^1$ -minimization to find (or closely approximate) the sparsest solution even when the conditions of these theorems do not hold. Whether or not this happens in the context of SRC is the topic of the next chapter.

## CHAPTER 11

### Equivalence on Highly-Coherent Data

In this chapter, we investigate whether sparsity is reliably achieved via  $\ell^1$ -minimization on highly-correlated data, such as class-structured databases.

#### 11.1. Inspiration

We are inspired by the data model and subsequent work of Wright and Ma [102] (see also the work of Wright et al. [103]), which produces an  $\ell^1/\ell^0$ -equivalence guarantee for dictionaries containing vectors assumed to model facial images. We summarize their result briefly.

Previous work has shown that the set of facial images of a fixed subject (person) under varying illumination conditions forms a convex cone, called an *illumination cone*, in pixel space [3, 44]. Wright and Ma demonstrate that in fact the set of facial images under varying illuminations over *all subjects combined* exhibits this cone structure. For example, they show that this is the case for the entire set of (raw) samples from the Extended Yale B Face Database [44]. Further, this cone becomes extremely narrow, i.e., a “bouquet,” as the number of pixels grows large [102]. These findings reiterate that class-structured data, particularly face databases, are highly-coherent.

Lee et al. [59] showed that any image from the illumination cone can be expressed as a linear combination of just a few images of the same subject under varying lighting conditions. In other words, illumination cones are well-approximated by linear subspaces. Thus the SRC condition that class manifolds are (approximately) linear subspaces presumably holds for databases made up of facial images under varying lighting conditions. Given a facial image  $\mathbf{y} \in \mathbb{R}^m$  that may be occluded or corrupted by noise,  $\mathbf{y}$  can thus be expressed as

$$(11.1) \quad \mathbf{y} = X_{\text{tr}}\boldsymbol{\alpha}_0 + \mathbf{z}_0,$$

given that certain requirements are satisfied in the sampling of the training data. By the above model,  $\boldsymbol{\alpha}_0$  is assumed to be non-negative (a result of the illumination cone model [44, 103]) and

## 11.1. INSPIRATION

---

sparse, containing nonzeros at training samples that represent the same subject as  $\mathbf{y}$  (i.e., are in the same class). Additionally,  $\mathbf{z}_0$  is an (unknown) error vector with nonzeros in only a fraction of its coordinates; i.e., the model assumes that only a portion of the pixels are occluded or corrupted [103].

The goal, as one might expect, is to recover  $\boldsymbol{\alpha}_0$  from Eq. (11.1). In the SRC paper [104], Wright et al. use  $\ell^1$ -minimization to do this. In particular, they solve

$$(11.2) \quad (\boldsymbol{\alpha}_1, \mathbf{z}_1) := \arg \min \|\boldsymbol{\alpha}\|_1 + \|\mathbf{z}\|_1 \text{ subject to } \mathbf{y} = X_{\text{tr}}\boldsymbol{\alpha} + \mathbf{z},$$

and they show that this “occlusion version” of SRC (previously mentioned in Section 5.4.3) produces very good classification results on corrupted facial images.

In a later paper, Wright, et al. [103] correctly note that the usual  $\ell^1/\ell^0$ -equivalence theorems do not hold on the highly-correlated data in  $X_{\text{tr}}$ , and so we are unable to determine whether or not the  $\ell^1$ -minimized solution  $\boldsymbol{\alpha}_1$  in Eq. (11.2) is equal to (what is assumed to be) the true sparsest solution  $\boldsymbol{\alpha}_0$ . Fortunately, Wright and Ma [102] proved a theorem that gives sufficient conditions for this equivalence under an assumed model (called the *bouquet model*) of facial images; see also Wright et al.’s version [103]. To state the theorem, we will need the following definition:

**DEFINITION 11.1.1** (Proportional Growth [102]). *A sequence of signal-error problems  $\mathbf{y} = X\boldsymbol{\alpha}_0 + \mathbf{z}_0$ , for  $X \in \mathbb{R}^{m \times N}$ , exhibits proportional growth with parameters  $\delta > 0$ ,  $\rho \in (0, 1)$ , and  $\beta > 0$ , if  $N = \lfloor \delta m \rfloor$ ,  $\|\mathbf{z}_0\|_0 = \lfloor \rho m \rfloor$ , and  $\|\boldsymbol{\alpha}_0\|_0 = \lfloor \beta m \rfloor$ .*

It follows that  $\delta$  is the redundancy factor in the dictionary  $X$  and  $\rho$  and  $\beta$  control the sparsity of  $\mathbf{z}_0$  and  $\boldsymbol{\alpha}_0$ , respectively. Here,  $\beta$  is assumed to be small and may depend on  $\delta$  and  $\rho$ .

We are now in a position to state Wright and Ma’s main theorem:

**THEOREM 11.1.1** (Wright and Ma [102], Theorem 1). *Fix any  $\delta > 0$  and  $\rho < 1$ . Suppose that  $X$  is distributed according to the bouquet model given by*

$$(11.3) \quad X = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{m \times N}, \quad \mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\boldsymbol{\mu}, (\nu^2/m)I_m), \quad \|\boldsymbol{\mu}\|_2 = 1, \quad \|\boldsymbol{\mu}\|_\infty \leq C_\mu m^{-1/2}, \quad C_\mu \geq 1$$

*for  $\nu$  sufficiently small. Also suppose that the sequence of signal-error problems  $\mathbf{y} = X\boldsymbol{\alpha}_0 + \mathbf{z}_0$  for  $X \in \mathbb{R}^{m \times N}$  exhibits proportional growth with parameters  $\delta$ ,  $\rho$ , and  $\beta$ . Suppose further that  $J \subset \{1, \dots, m\}$  is a uniform random subset of size  $\rho m$ , and that  $\boldsymbol{\sigma} \in \mathbb{R}^m$  with entries of  $\boldsymbol{\sigma}_J$*

## 11.2. PROJECT DESCRIPTION

---

*i.i.d.*  $\pm 1$  (*independent of J*) and  $\boldsymbol{\sigma}_{JC} = \mathbf{0}$ . Lastly assume that  $m$  is sufficiently large. Then with probability at least  $1 - C \exp(-\zeta^* m)$  in  $X$ ,  $J$ , and  $\boldsymbol{\sigma}$ , for all  $\boldsymbol{\alpha}_0$  with  $\|\boldsymbol{\alpha}_0\|_0 \leq \beta^* m$  and any  $\mathbf{z}_0$  with sign vector  $\boldsymbol{\sigma}$  and support  $J$ , we have

$$(\boldsymbol{\alpha}_0, \mathbf{z}_0) = \arg \min_{\boldsymbol{\alpha}, \mathbf{z}} \|\boldsymbol{\alpha}\|_1 + \|\mathbf{z}\|_1 \text{ subject to } X\boldsymbol{\alpha} + \mathbf{z} = X\boldsymbol{\alpha}_0 + \mathbf{z}_0.$$

Here,  $C$  is a numerical constant and  $\beta^*$  and  $\zeta^*$  are positive constants (independent of  $m$ ) which depend on  $\delta$ ,  $\rho$ , and  $\nu$ . By “ $\nu$  sufficiently small” and “ $m$  sufficiently large,” Wright and Ma mean that there exist constants  $0 < \nu < \nu^*$  and  $m > m^*$  (independent of  $m$ ) such that  $\nu^*(\delta, \rho) > 0$  and  $m^*(\delta, \rho, \nu) > 0$ , respectively.<sup>1</sup>

This theorem illustrates that  $\ell^1/\ell^0$ -equivalence can provably hold on highly-coherent data!

### 11.2. Project Description

Despite its applicability to highly-coherent data, Theorem 11.1.1 does not prove that  $\ell^1/\ell^0$ -equivalence holds in SRC. First of all, the theorem requires that  $m$  be sufficiently large, which may not be the case when dimensionality reduction is used. Second, the model in Theorem 11.1.1 does not explicitly deal with class-structured data. A true face recognition model should account for the individual subjects, with samples in the same class being (on average) more correlated than those from different classes. Thus our model should contain “sub-bouquets” (i.e., the classes) inside the larger bouquet.

With these changes in mind, we design an experiment to study the relationship between sparsity and  $\ell^1$ -minimization on highly-coherent and class-structured data, such as the images used in face recognition. First, we specify the dimension  $m$ , the number of samples in each training class  $N_0$ , and the number of classes  $L$ . We require that  $N_{\text{tr}} := N_0 L > m$  so that the resulting dictionary of training samples leads to an underdetermined system. We then randomly generate training data with an increasing amount of cone/bouquet structure as well as class structure, along with a test sample—with known sparse coefficient vector  $\boldsymbol{\alpha}_0$ —generated as a linear combination of training samples from a single class. We run a fixed number of trials of the experiment at each of 11

---

<sup>1</sup>The relationship between  $\beta^*$  and  $\beta$  is not explicitly stated, but it makes sense that  $\beta^* \leq \beta$  by the proportional growth assumption. Further, if  $\beta = \beta(\delta, \rho)$ , then since  $\beta^* = \beta^*(\delta, \rho, \nu)$ , we can likely alternatively write  $\beta^* = \beta^*(\beta, \nu)$ .

### 11.3. EXPERIMENTS

---

increasing values of coherence (we call these *stages*) and determine at which stages  $\ell^1$ -minimization can closely (or exactly) recover  $\boldsymbol{\alpha}_0$ .

#### 11.3. Experiments

**11.3.1. Experimental Setup.** For each generated training set  $X_{\text{tr}} = [X^{(1)}, \dots, X^{(L)}] \in \mathbb{R}^{m \times N_{\text{tr}}}$ , we set the (clean) test sample  $\mathbf{y}_0$  to be a random vector in the positive span of the class 1 data. That is, we set

$$\mathbf{y}_0 := \alpha_1^{(1)} \mathbf{x}_1^{(1)} + \dots + \alpha_{N_0}^{(1)} \mathbf{x}_{N_0}^{(1)},$$

where  $X^{(1)} := [\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_{N_0}^{(1)}]$  and  $\alpha_j^{(1)} \sim \text{unif}(0, 1)$ ,  $1 \leq j \leq N_0$ . We then define

$$\boldsymbol{\alpha}_0 := [\alpha_1^{(1)}, \dots, \alpha_{N_0}^{(1)}, 0, \dots, 0]^T \in \mathbb{R}^{N_{\text{tr}}}.$$

Given this setup, we want to see if  $\ell^1$ -minimization will recover  $\boldsymbol{\alpha}_0$ , i.e., if the solution

$$\boldsymbol{\alpha}_1 := \arg \min_{\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}} \|\boldsymbol{\alpha}\|_1 \text{ subject to } X_{\text{tr}} \boldsymbol{\alpha} = \mathbf{y}_0$$

is equal to  $\boldsymbol{\alpha}_0$ . Note that for large  $L$ ,  $\boldsymbol{\alpha}_0$  can be viewed as a sparse vector.

In Stage 1 of our model, the training data has no class or cone structure and is randomly generated on the unit sphere  $S^{m-1}$ . It has been shown experimentally that, for  $N_{\text{tr}} = 2m$  and  $m$  sufficiently large, an  $\ell^1$ -minimization solution with no more than  $(3/10)m$  nonzeros is enough to ensure it is the sparsest solution with high probability [29]. Thus we expect to see exact recovery in Stage 1 for values of  $N_0$ ,  $m$ , and  $L$  satisfying these requirements.

To add both bouquet and class (or sub-bouquet) structure to the training set in subsequent stages, we define the cone mean  $\bar{\mathbf{x}}$  and the class means  $\{\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_L\}$ . At Stage  $i$ ,  $1 \leq i \leq 11$ , we set  $\bar{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, I_m)$  and then modify  $\bar{\mathbf{x}} \leftarrow \mu_i \bar{\mathbf{x}} / \|\bar{\mathbf{x}}\|_2$ , where  $\mu_i := (i-1)/10$  effectively increases the cone mean from 0 as  $i$  increases. Next, each class mean is randomly generated depending on  $\bar{\mathbf{x}}$  as follows: For each class  $1, \dots, L$ , we sample  $\bar{\mathbf{x}}_l$  from  $\mathcal{N}(\bar{\mathbf{x}}, \eta_l m^{-1/2} I_m)$  for  $\eta_l := 2/l$  (so that each class mean becomes increasingly close to the cone mean) and then modify  $\bar{\mathbf{x}}_l \leftarrow \mu_l \bar{\mathbf{x}}_l / \|\bar{\mathbf{x}}_l\|_2$ ,  $1 \leq l \leq L$ . Lastly, to generate the training samples in class  $1 \leq l \leq L$ , we sample  $\mathbf{x}_j^{(l)}$  from

### 11.3. EXPERIMENTS

---

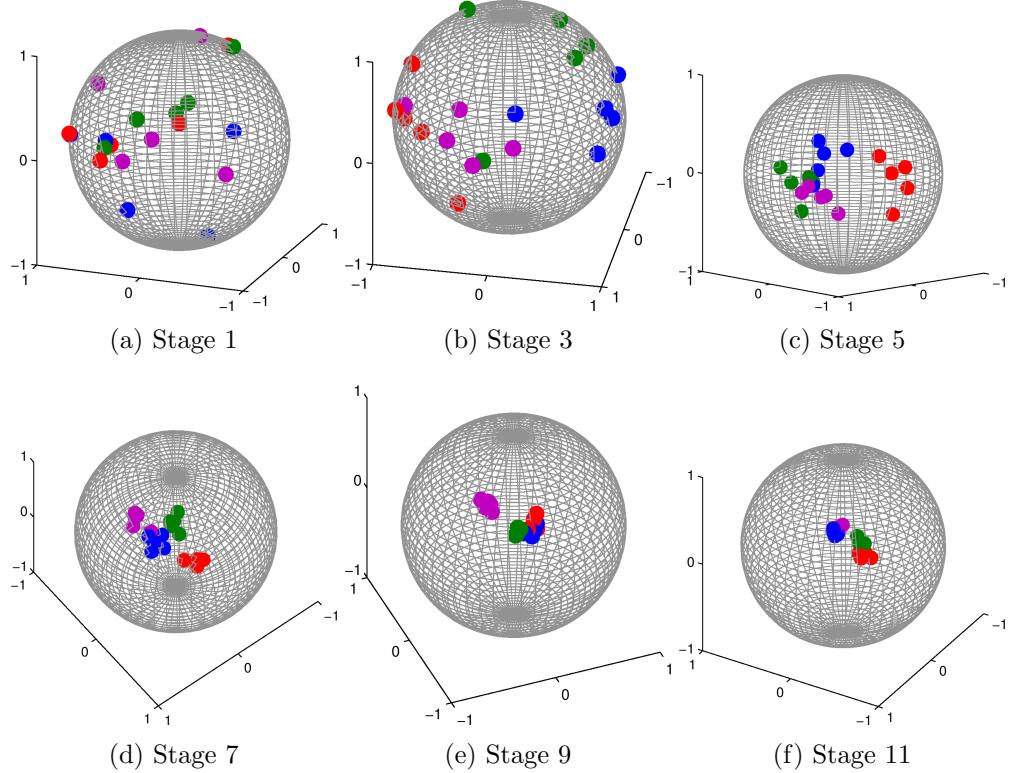


FIGURE 11.1. An example of the generated training data from the random database model across odd-numbered stages (as mutual coherence increases) with  $m = 3$ ,  $N_0 = 5$ , and  $L = 4$ . The colors denote the classes. Plots have been manually rotated to aid in visualization. (a) At Stage 1, data is uniformly spread out on the sphere; (b)-(f) At increasingly higher stages, the data set as a whole becomes more bouquet-shaped, as does the data in each class.

$\mathcal{N}(\bar{\mathbf{x}}_l, (\eta_i m^{-1/2}/L) I_m)$  and then modify  $\mathbf{x}_j^{(l)} \leftarrow \mathbf{x}_j^{(l)} / \|\mathbf{x}_j^{(l)}\|_2$ ,  $1 \leq j \leq N_0$ . Figure 11.1 shows an example of Stage  $i \in \{1, 3, \dots, 11\}$  with  $m = 3$ ,  $N_0 = 5$ , and  $L = 4$ .

We perform experiments using four different specifications for the triples  $(N_0, m, L)$ , as shown in Table 11.1. By design, we have that  $\|\boldsymbol{\alpha}_0\|_0 = N_0$  in our experiments (though we will also briefly look at the case that  $\|\boldsymbol{\alpha}_0\|_0 < N_0$ ). Note that (i) the inequality  $\|\boldsymbol{\alpha}_0\|_0 < (3/10)m$  is satisfied for each of the specifications in Table 11.1 and (ii) these numbers are similar to what we might expect to see in classification of a face database (after some method of feature extraction is applied, as is generally required by SRC for face classification).

#### 11.3.2. Experimental Results: No Noise.

### 11.3. EXPERIMENTS

---

ID	$(N_0, m, L)$	$\ \boldsymbol{\alpha}_0\ _0/m$	$\ \boldsymbol{\alpha}_0\ _0/N_{\text{tr}}$	Redundancy	Comments
DB-1	(5,50,20)	1/10	1/20	2:1	Baseline redundancy; $N_0$ small with respect to $m$ , $N_{\text{tr}}$
DB-2	(10,50,10)	1/5	1/10	2:1	Baseline redundancy; $N_0$ less small with respect to $m$ , $N_{\text{tr}}$
DB-3	(10,50,50)	1/5	1/50	10:1	Highly redundant; large $L$
DB-4	(5,200,50)	1/40	1/50	5:4	Barely redundant; large $L$

TABLE 11.1. Specification of parameters in the random database model.

*Accuracy of recovery.* We consider the following quantities for evaluating the success of  $\ell^1/\ell^0$ -recovery:

- The average normalized  $\ell^2$ -error  $\text{err}_{\ell^2} := \|\boldsymbol{\alpha}_1 - \boldsymbol{\alpha}_0\|_2 / \|\boldsymbol{\alpha}_0\|_2$  between the  $\ell^1$ -minimized solution  $\boldsymbol{\alpha}_1$  and  $\boldsymbol{\alpha}_0$ ,
- The average number of nonzeros of  $\boldsymbol{\alpha}_1$  occurring at training samples *not* in class 1 (we call these “off-support” nonzeros, because they are nonzeros not in the support of  $\boldsymbol{\alpha}_0$ ), divided by the total number of nonzeros. That is, let  $\boldsymbol{\alpha}_1^{\text{off-supp}}$  be the result of setting all entries in  $\boldsymbol{\alpha}_1$  that are in class 1 to zero. Then this error is defined as

$$\text{err}_{\text{supp}} := \frac{\|\boldsymbol{\alpha}_1^{\text{off-supp}}\|_0}{\|\boldsymbol{\alpha}_1\|_0},$$

- Since  $\text{err}_{\text{supp}}$  does not provide information regarding the *size* of the off-support nonzero coefficients, we also consider

$$\text{err}_{\text{supp}(\ell^2)} := \frac{\|\boldsymbol{\alpha}_1^{\text{off-supp}}\|_2}{\|\boldsymbol{\alpha}_1\|_2} \quad \text{and} \quad \text{err}_{\text{supp}(\ell^1)} := \frac{\|\boldsymbol{\alpha}_1^{\text{off-supp}}\|_1}{\|\boldsymbol{\alpha}_1\|_1},$$

- The average mutual coherence of the training set,  $\mu(X_{\text{tr}}) =: \mu$ .

It is informative to consider the effect that the support error quantities would (hypothetically) have on the classification performance of SRC. Recall that SRC computes the class residuals  $\text{err}_l(\mathbf{y}_0) := \|\mathbf{y}_0 - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}_1)\|_2$  (in the case that the clean test sample  $\mathbf{y}_0$  is known) and assigns  $\mathbf{y}_0$  to the class with the smallest residual. Thus if  $\text{err}_{\text{supp}}$ ,  $\text{err}_{\text{supp}(\ell^2)}$ , and  $\text{err}_{\text{supp}(\ell^1)}$  are small, we expect that SRC will have an easier time classifying the test sample correctly. For example, if all the support error quantities are 0, then  $\delta_1(\boldsymbol{\alpha}_1) = \boldsymbol{\alpha}_1$  and it follows that the class 1 residual  $\text{err}_1(\mathbf{y}_0) = 0$  and  $\text{err}_l(\mathbf{y}_0) = \|\mathbf{y}_0\|_2$  for  $2 \leq l \leq L$ . This corresponds to the ideal classification scenario.

We compute the average quantities  $\text{err}_{\ell^2}$ ,  $\text{err}_{\text{supp}}$ ,  $\text{err}_{\text{supp}(\ell^2)}$ ,  $\text{err}_{\text{supp}(\ell^1)}$ , and  $\mu$  over 1000 trials at each stage, using the  $\ell^1$ -minimization algorithm HOMOTOPY with error/sparsity trade-off

### 11.3. EXPERIMENTS

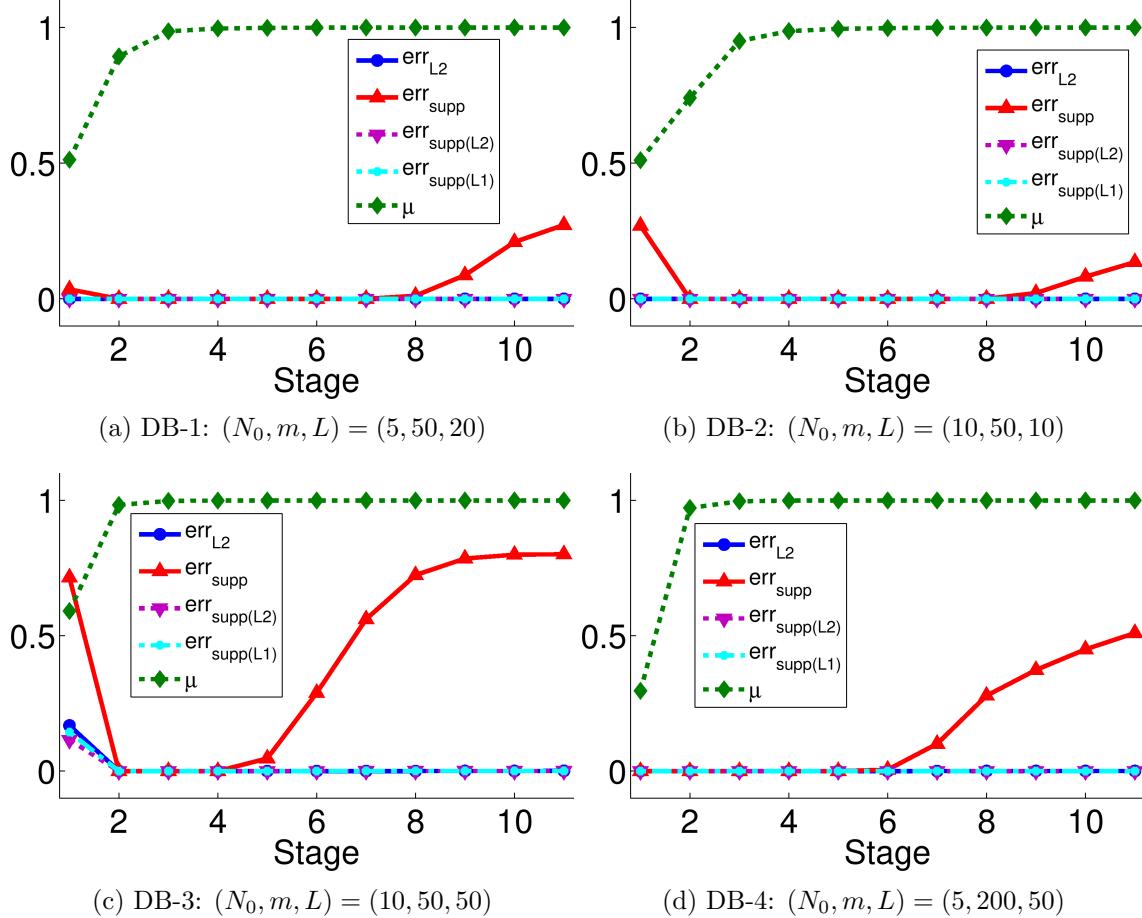


FIGURE 11.2. Recovery results on random database model (average of 1000 trials) in the case of no noise.

parameter  $\lambda = 10^{-10}$  (to force near-exactness in the approximation). The results are shown in Figure 11.2.

Considering that  $\text{err}_{\text{supp}}$  records *any* off-support nonzeros, regardless of how small, the results are quite good. In many cases,  $\ell^1$ -minimization was able to recover the exact solution  $\alpha_0$  on highly-correlated data, and when errors in the support occurred, they were generally small.

We see two different things happening at either end of the Stage axis. At Stage 1, we see support errors in every database except DB-4 (the low-redundancy case). Further, there are nonzero values of  $\text{err}_{\ell^2}$ ,  $\text{err}_{\text{supp}(\ell^2)}$ , and  $\text{err}_{\text{supp}(\ell^1)}$  for DB-3 (the high-redundancy case) at this stage. At high stages, we see similar small support errors as the data became very correlated; these support errors were numerous (accounting for around half the nonzero coefficients) for both DB-3 and DB-4.

### 11.3. EXPERIMENTS

---

We start by explaining the results at Stage 1. Given the plots in Figure 11.2, our instinct may be to suspect that something wrong happened here, especially considering the exact recovery on all databases at Stage 2. For the cases that we had a ratio of 2-to-1 redundancy, does this contradict the experimental result [29] that having  $N_0 = \|\boldsymbol{\alpha}_0\|_0 < (3/10)m$  nonzeros guarantees  $\ell^1/\ell^0$ -equivalence with high probability? It would, but for the fact that this result holds asymptotically. To test this, we repeated the experiments for increasing values of  $m$ , scaling  $N_0$  and  $L$  accordingly so that the redundancy remained constant. More precisely, we defined  $r_1 := m/N_{\text{tr}}$  and  $r_2 := N_0/L$  and then set  $\tilde{L} := [\sqrt{\tilde{m}/(r_1 r_2)}]$  and  $\tilde{N}_0 := r_2 \tilde{L}$ . Here,  $[\cdot]$  denotes the nearest integer function and  $\tilde{m}$ ,  $\tilde{L}$ , and  $\tilde{N}_0$  denote the increased values of  $m$ ,  $L$ , and  $N_0$ , respectively. As we illustrate in Figure 11.3, the value of  $\text{err}_{\text{supp}}$  decreased to 0 as  $\tilde{m}$  increased. As is to be expected, both the amount of redundancy and the relationship  $N_0/m$  affected the speed of convergence. We exclude results for DB-4, as we already see perfect recovery at Stage 1 in Figure 11.2d.

In comparing the Stage 1 results to those from data with bouquet/cone structure (i.e., Stages 2-11), it is initially surprising that small to moderate levels of correlation in the data samples appear to *improve* sparse recovery. As mentioned, we see near-perfect recovery of  $\boldsymbol{\alpha}_0$  at Stage 2 for every tried  $(N_0, m, L)$  triple; this is in stark contrast to the recovery accuracy at Stage 1, especially for DB-3 (Figure 11.2c). This sharp change coincides with a significant increase in the within-class correlation between Stages 1 and 2 in our model, whereas the correlation between classes essentially remains unchanged. Though the exact specifics will depend on the  $\ell^1$ -minimization algorithm used, we strongly suspect that the relative clustering of the samples in the support of  $\boldsymbol{\alpha}_0$  at Stage 1 (as compared to their random distribution at Stage 1) make it much easier for the algorithm to recover the desired solution.

Conversely, at high stages, it appears that the loss of class structure negatively affected the recovery of  $\boldsymbol{\alpha}_0$ . As the standard deviation of the class mean distributions grew small, the class cones began to significantly overlap, and  $\ell^1$ -minimization could not exactly recover the support of  $\boldsymbol{\alpha}_0$ . Notice that we see an especially large number of support errors  $\text{err}_{\text{supp}}$  for databases with large values of  $L$ , namely, DB-3 (Figure 11.2c) and DB-4 (Figure 11.2d). For DB-3, the nonzero values of  $\text{err}_{\text{supp}}$  at Stages 5 and 6 (compared to  $\text{err}_{\text{supp}} \approx 0$  at these stages for DB-4) confirms that redundancy, as well as the number of classes, affects recovery.

### 11.3. EXPERIMENTS

---

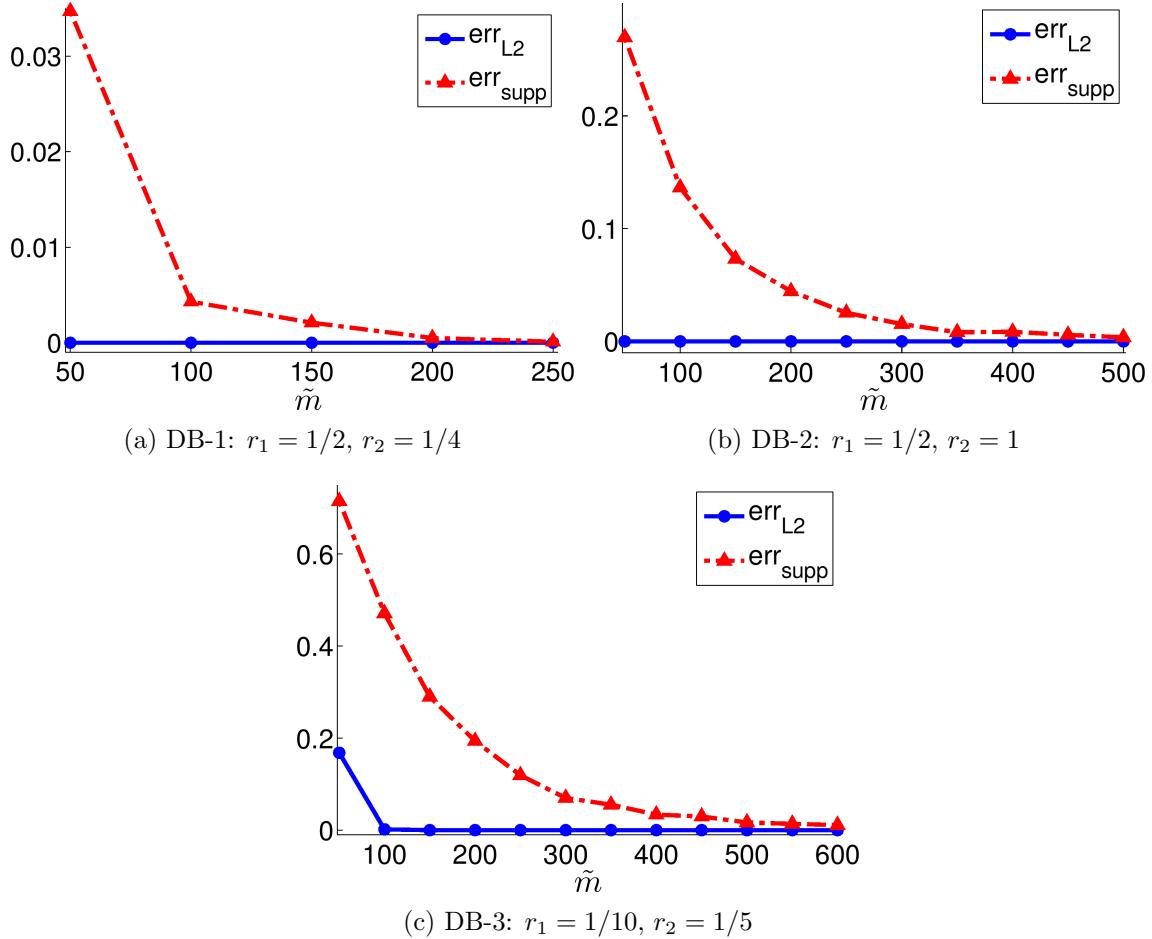


FIGURE 11.3. Asymptotic recovery at Stage 1 of the random database model (average of 1000 trials). Note the different scales.

*Effect on classification.* We earlier discussed the relationship between the support error quantities  $\text{err}_{\text{supp}}$ ,  $\text{err}_{\text{supp}(\ell^2)}$ , and  $\text{err}_{\text{supp}(\ell^1)}$  on the classification performance of SRC, in particular, their effect on the class residuals  $\text{err}_l(\mathbf{y}_0) := \|\mathbf{y}_0 - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}_1)\|_2$ . Here, we consider these residuals explicitly. For each of the four databases, we computed the average residual  $\text{err}_l(\mathbf{y}_0)$  (over 1000 trials) for each class  $1 \leq l \leq L$  at each of the 11 values of coherence.

Not surprisingly given the small support error quantities determined in the previous section, there is a stark difference between the residual of class 1 and those of the other classes at all stages. More precisely, the ideal classification scenario occurs in all cases, with  $\text{err}_1(\mathbf{y}_0) \approx 0$  and  $\text{err}_l(\mathbf{y}_0) \approx \|\mathbf{y}_0\|_2$  for all  $2 \leq l \leq L$ . The approximations are of the order  $10^{-8}$  (or better), except

### 11.3. EXPERIMENTS

---

for the highly-redundant database DB-3 at Stage 1. In this case, the average quantities were  $\text{err}_1(\mathbf{y}_0) = 0.230$  and

$$\|\mathbf{y}_0\|_2 - \underset{2 \leq l \leq L}{\text{mean}} \text{err}_l(\mathbf{y}_0) = 0.004.$$

These findings are consistent with the results in Figure 11.2c. Even though these quantities at Stage 1 are nonzero, it is important to note that good classification would still be achieved, as  $\min_{2 \leq l \leq L} \text{err}_l(\mathbf{y}_0) = 1.806$ , which is a lot larger than  $\text{err}_1(\mathbf{y}_0) = 0.230$ .

*Varying the sparsity level.* We next consider what happens when the sparsity level  $\|\boldsymbol{\alpha}_0\|_0$  is strictly less than the number of class 1 training samples  $N_0$ . This is important to investigate: can  $\ell^1$ -minimization identify the correct training samples from among the rest of the (highly-correlated) training data in that class? For DB-2 and DB-3, we generated  $\boldsymbol{\alpha}_0$  (and subsequently  $\mathbf{y}_0$ ) using the first five samples in class 1. Figures 11.4a and 11.4c show the recovery results, and Figures 11.4b and 11.4d repeat the plots in Figures 11.2b and 11.2c (in which  $\|\boldsymbol{\alpha}_0\|_0 = N_0$ ) for convenient comparison.

At Stage 1, we see that the support of  $\boldsymbol{\alpha}_1$  was more concentrated on the correct training samples when  $\|\boldsymbol{\alpha}_0\|_0$  was smaller, evidenced by smaller values of  $\text{err}_{\text{supp}}$ . This is to be expected, as the ground truth solution became sparser. For the lower-redundancy case DB-2, we see far more support errors as the correlation increased when  $\|\boldsymbol{\alpha}_0\|_0 = 5$  (Figure 11.4a) than for the case  $\|\boldsymbol{\alpha}_0\|_0 = N_0$  (Figure 11.4b); however, the values of these off-support coefficients were very small, as demonstrated by the near-zero values of  $\text{err}_{\text{supp}(\ell^2)}$  and  $\text{err}_{\text{supp}(\ell^1)}$ . Though class 1 training samples not in the support of  $\boldsymbol{\alpha}_0$  were mistakenly selected as the data in class 1 became more correlated, these samples played a negligible role in the representation. For the high-redundancy case DB-3, we similarly see more small-valued, off-support coefficients at Stages 2–5 when  $\|\boldsymbol{\alpha}_0\|_0 = 5$  (Figure 11.4c) than in the case  $\|\boldsymbol{\alpha}_0\|_0 = N_0$  (Figure 11.4d). The value of  $\text{err}_{\text{supp}}$  for  $\|\boldsymbol{\alpha}_0\|_0 < N_0$  was actually smaller than it was for  $\|\boldsymbol{\alpha}_0\|_0 = N_0$  at many of the higher stages, however, suggesting that the added degree of sparsity helped to counter-balance the high redundancy of this database (and its negative effect on recovery) in these cases.

*Eliminating errors by thresholding.* Before we turn to the noisy setting, we demonstrate that the small support errors in  $\boldsymbol{\alpha}_1$  depicted in Figure 11.2 can be completely remedied using thresholding in

### 11.3. EXPERIMENTS

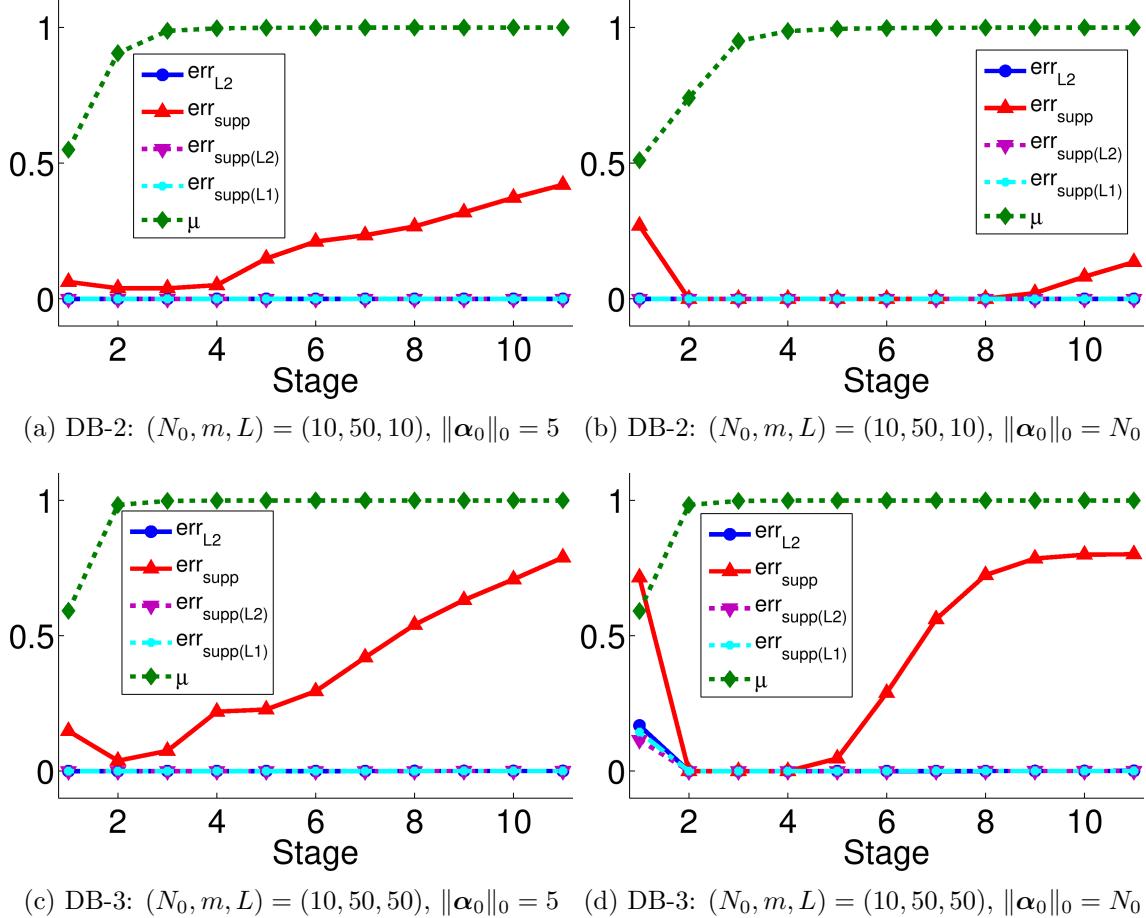


FIGURE 11.4. Comparing  $\|\boldsymbol{\alpha}_0\|_0 < N_0$  and  $\|\boldsymbol{\alpha}_0\|_0 = N_0$  sparsity levels (average of 1000 trials) on the random database model in the case of no noise.

all but the high-redundancy case DB-3. After determining  $\boldsymbol{\alpha}_1$  as before, we set its small coefficients (those with absolute value less than some threshold  $\tau$ ) to zero, obtaining the vector  $\boldsymbol{\alpha}_1^\tau$ . We then re-solved the equation  $X_{\text{tr}}\boldsymbol{\alpha} = \mathbf{y}_0$  with the constraint that the solution, denoted  $\hat{\boldsymbol{\alpha}}_1$ , had the same support as the thresholded  $\boldsymbol{\alpha}_1^\tau$ . For simplicity, we did this by setting the columns of  $X_{\text{tr}}$  corresponding to zero-coordinates in  $\boldsymbol{\alpha}_1^\tau$  to  $\mathbf{0}$ , thus obtaining the matrix  $\hat{X}_{\text{tr}}$ . We then used MATLAB's “\” operator to define  $\hat{\boldsymbol{\alpha}}_1 := \hat{X}_{\text{tr}} \setminus \mathbf{y}_0$ . In our case, since  $X_{\text{tr}}$  was not square, the desired least squares solution was found by (MATLAB's implementation of) QR-factorization.

For all but the highly-redundant database DB-3,  $\hat{\boldsymbol{\alpha}}_1$  was equal to the sparsest solution  $\boldsymbol{\alpha}_0$  (up to nearly machine-precision) for the thresholding value  $\tau = 10^{-5}$ . For  $\tau \in \{0.001, 0.01\}$  on these three databases (DB-1, DB-2, and DB-3), we saw small nonzero values of  $\text{err}_{\ell^2}$ , but these errors

### 11.3. EXPERIMENTS

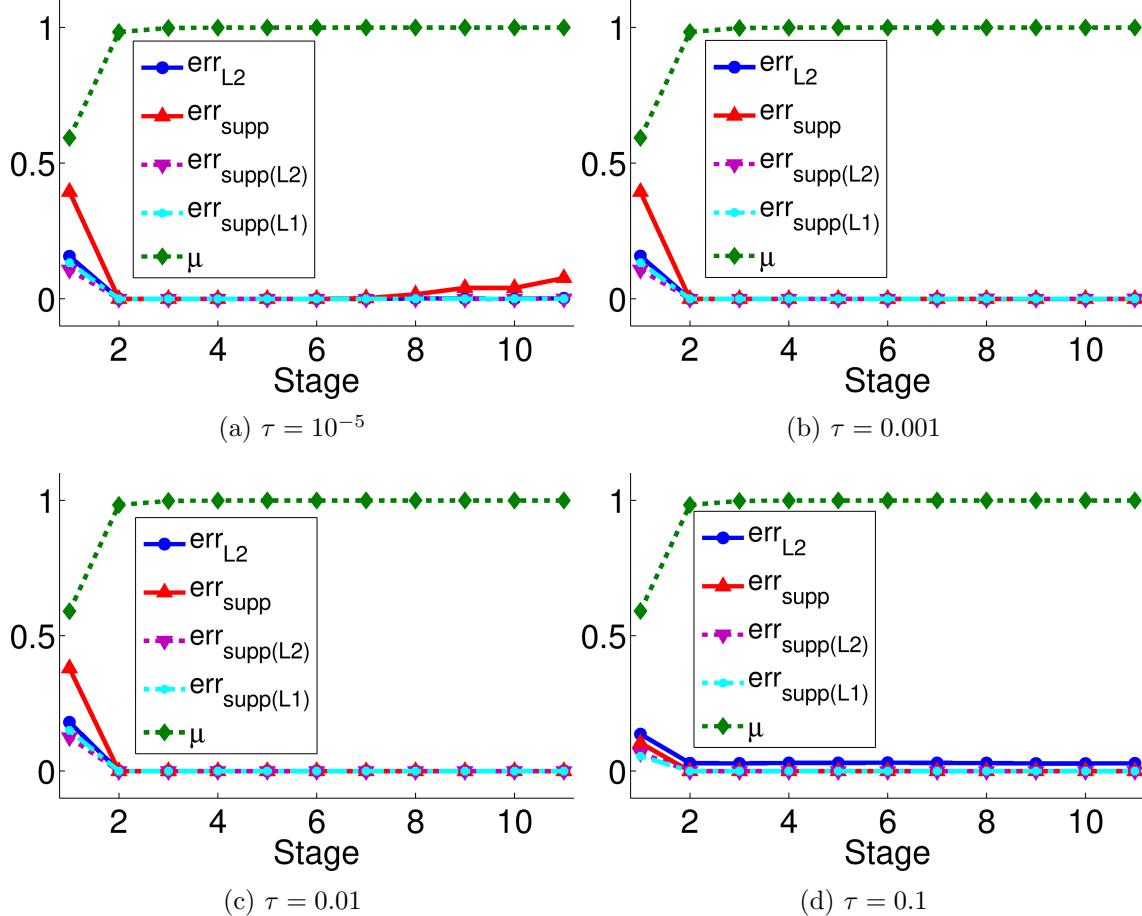


FIGURE 11.5. The results of thresholding (average of 1000 trials) on the highly-redundant database DB-3:  $(N_0, m, L) = (10, 50, 50)$  in the case of no noise.

were indiscernible in plots on the same scale as those in Figure 11.2, and so we do not show them here. For  $\tau = 0.1$ , there was a consistent, small but nontrivial  $\ell^2$ -error across all stages, as small coefficients corresponding to class 1 training samples were incorrectly set to 0. For all four values of  $\tau$ , there were no support errors.

For the high-redundancy case DB-3, we continued to see errors at Stage 1, similar to those in Figure 11.2c. For the thresholding values  $\tau \in \{0.001, 0.01, 0.1\}$  (i.e., for  $\tau$  large enough), there were no support errors at other stages. However, similarly to the other databases, we saw nontrivial  $\ell^2$ -error when  $\tau = 0.1$ . We plot the results for DB-3 in Figure 11.5, stressing that the results for the other databases contained errors too small to produce nontrivial plots.

### 11.3. EXPERIMENTS

---

**11.3.3. Experimental Results: Noisy Setting.** In these experiments, we examine  $\ell^1/\ell^0$ -recovery when noise is added to the test sample  $\mathbf{y}_0$ . Recall the theorems by Donoho et al. regarding  $\ell^1/\ell^0$ -equivalence in the noisy setting stated in Theorems 9.2.5 and 9.2.6.

*Accuracy of recovery.* Unfortunately, Eq. (9.6) and Eq. (9.7) in the referenced theorems do not make sense for large mutual coherence  $\mu(X)$ . However, we can still look for a correlation between  $\|\boldsymbol{\alpha}_0 - \boldsymbol{\alpha}_{1,\epsilon}\|_2$  (where  $\boldsymbol{\alpha}_{1,\epsilon}$  is the solution to Eq. (11.4) below) and the values of the noise tolerance  $\zeta$  (see that statement of Theorem 9.2.5), the approximation error bound  $\epsilon$ ,  $N_0 = \|\boldsymbol{\alpha}_0\|_0 := k$ , and  $\mu(X_{\text{tr}})$ , with  $\epsilon =: C\zeta$  for some constant  $C > 0$ . We modify the experiments in Section 11.3.2 as follows: First, we specify the noise tolerance  $\zeta$  and the constant  $C$ . After generating the training data and the (noise-free) test sample  $\mathbf{y}_0$ , we set  $\mathbf{y} := \mathbf{y}_0 + \mathbf{z}$ , where the entries of  $\mathbf{z}$  are drawn from  $\mathcal{N}(0, \zeta/(2\sqrt{m}))$ . Then  $\|\mathbf{z}\|_2 \leq \zeta$  with probability at least 95%. From here, we set  $\epsilon := C\zeta$  and find

$$(11.4) \quad \boldsymbol{\alpha}_{1,\epsilon} := \arg \min_{\boldsymbol{\alpha}} \|\boldsymbol{\alpha}\|_1 \text{ subject to } \|\mathbf{y} - X_{\text{tr}}\boldsymbol{\alpha}\|_2 \leq \epsilon.$$

We set  $\zeta = 0.01$ , and we used two values of  $C$ :  $C = 5$ , and  $C = 10$ , producing the  $(\zeta, \epsilon)$ -pairs  $(0.01, 0.05)$  and  $(0.01, 0.1)$ . In order to ensure that the reconstruction error  $\|X_{\text{tr}}\boldsymbol{\alpha}_{1,\epsilon} - \mathbf{y}\|_2$  was less than  $\epsilon$ , we used the basis pursuit denoising version of the  $\ell^1$ -minimization algorithm SPGL1 [95, 96].

In Figure 11.6, we plot the normalized  $\ell^2$ -error, the fraction of off-support nonzeros, the normalized  $\ell^2$  and  $\ell^1$ -norms of the off-class support vectors, and the mutual coherence  $\mu(X_{\text{tr}}) =: \mu$ . Note that we modify the corresponding definitions given in Section 11.6 (for  $\text{err}_{\ell^2}$ ,  $\text{err}_{\text{supp}}$ ,  $\text{err}_{\text{supp}(\ell^2)}$  and  $\text{err}_{\text{supp}(\ell^1)}$ ) to use  $\boldsymbol{\alpha}_{1,\epsilon}$  instead of  $\boldsymbol{\alpha}_1$  and do not change the notation. We report the averages over 1000 trials at each stage.

As we can see, there is clearly a relationship between  $\text{err}_{\ell^2}$  and the amount of correlation in the data. As the data became increasingly bouquet-shaped, both within each class and as a data set as a whole, the normalized  $\ell^2$ -distance between  $\boldsymbol{\alpha}_{1,\epsilon}$  and  $\boldsymbol{\alpha}_0$  increased. The rate of increase of this error appears to be related the redundancy of the database. It is evident that mutual coherence was not a good indicator of  $\text{err}_{\ell^2}$ , as the plots show that  $\text{err}_{\ell^2}$  could be relatively low even after  $\mu(X_{\text{tr}})$  had reached its maximum value.

Perhaps more importantly, the supports of the solution vectors  $\boldsymbol{\alpha}_{1,\epsilon}$  and  $\boldsymbol{\alpha}_0$  were nearly identical at stages greater than 1. This means that the vast majority of nonzeros in  $\boldsymbol{\alpha}_0$  occurred at positions

### 11.3. EXPERIMENTS

---

corresponding to class 1 training samples. To fix the small support errors, we could use the thresholding technique discussed in the previous section, choosing  $\tau$  by trial-and-error. This method could also be used to ameliorate the numerous support errors for the databases DB-2 and DB-3 at Stage 1. In this case, we found that  $\tau = 0.01$  greatly reduced the Stage 1 support errors but did not eliminate them completely.

Lastly, we consider the differences between setting  $C = 5$  and  $C = 10$ . For the most part, the plots are quite similar. We see that setting  $C = 5$  produced slightly better recovery than  $C = 10$  at Stage 1, but in general, the normalized  $\ell^2$ -error  $\text{err}_{\ell^2}$  was the same for the two settings at higher stages. This is very informative, as it tells us that  $\ell^1/\ell^0$ -recovery on this kind of highly-correlated data is potentially quite robust to the setting of  $C$  in the approximation error tolerance  $\epsilon = C\zeta$ . Once again, we attribute this to the class structure of the data making it easier for the  $\ell^1$ -minimization algorithm to find the class solution  $\boldsymbol{\alpha}_0$ .

*Effect on classification.* As in the noise-free scenario, we compute the class residuals  $\text{err}_l(\mathbf{y}) := \|\mathbf{y} - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}_{1,\epsilon})\|_2$  for each of the four databases at each of the 11 values of coherence. Specifically, we are interested in how close the class 1 residual is to 0 (signifying perfect reconstruction of  $\mathbf{y}$  using class 1) and how close the next smallest class residual  $\min_{2 \leq l \leq L} \text{err}_l(\mathbf{y})$  is to this value. If it is close, then it means that we should have less confidence in the SRC classification assignment than if these quantities were far apart, i.e., that SRC distinguishes the correct class less clearly.

The average relevant class residuals (over 1000 trials) are displayed in Table 11.2. Since the results for  $C = 5$  and  $C = 10$  were very similar, we only include the results for  $C = 5$ .

Noting that  $\epsilon := C\zeta = 0.05$ , we see that the ideal classification scenario occurred in nearly all cases. That is, since  $\text{err}_1(\mathbf{y}) \approx \epsilon$  almost always, class 1 training samples made up essentially the entire approximation of the test sample. The exception, again, was DB-3 at Stage 1, for which  $\text{err}_1(\mathbf{y})$  and  $\min_{2 \leq l \leq L} \text{err}_l(\mathbf{y})$  were the least separated (i.e., relatively close in value). However, correct classification would still be achieved.

The reader might notice that the quantities  $\min_{2 \leq l \leq L} \text{err}_l(\mathbf{y})$  at Stage 1 are lower than at higher stages; this is because

$$\min_{2 \leq l \leq L} \text{err}_l(\mathbf{y}) = \min_{2 \leq l \leq L} \|\mathbf{y} - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}_{1,\epsilon})\|_2 \approx \|\mathbf{y} - X_{\text{tr}}\mathbf{0}\|_2 = \|\mathbf{y}\|_2$$

### 11.3. EXPERIMENTS

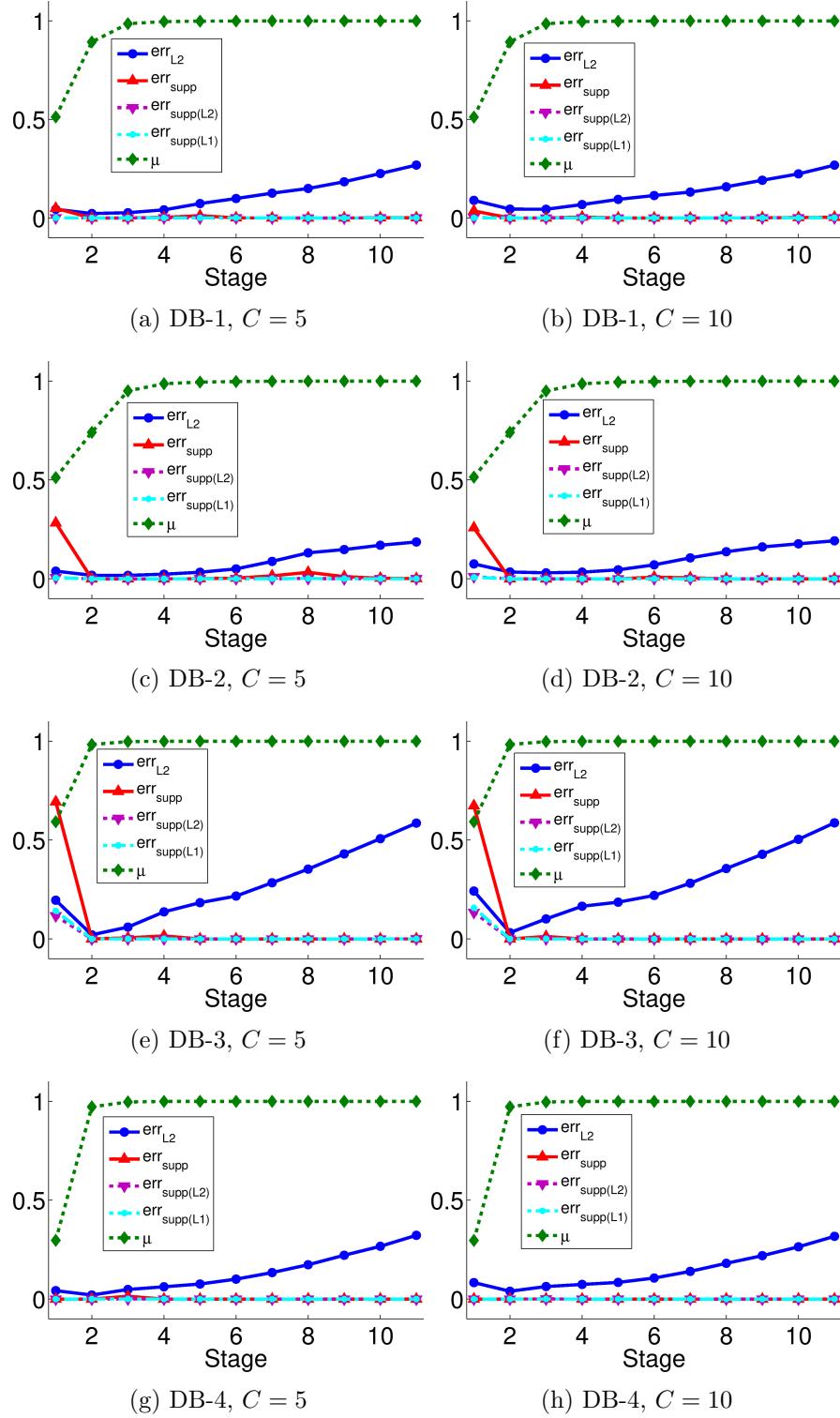


FIGURE 11.6. Recovery results on the random database model in the case of noise.

## 11.4. SUMMARY

---

Stage	DB-1		DB-2		DB-3		DB-4	
	err <sub>1</sub> ( $\mathbf{y}$ )	min <sub>2 ≤ l ≤ L</sub> err <sub>l</sub> ( $\mathbf{y}$ )	err <sub>1</sub> ( $\mathbf{y}$ )	min <sub>2 ≤ l ≤ L</sub> err <sub>l</sub> ( $\mathbf{y}$ )	err <sub>1</sub> ( $\mathbf{y}$ )	min <sub>2 ≤ l ≤ L</sub> err <sub>l</sub> ( $\mathbf{y}$ )	err <sub>1</sub> ( $\mathbf{y}$ )	min <sub>2 ≤ l ≤ L</sub> err <sub>l</sub> ( $\mathbf{y}$ )
1	0.05	1.27	0.06	1.81	0.28	1.80	0.05	1.27
2	0.05	2.30	0.05	3.77	0.05	4.92	0.05	2.47
3	0.05	2.47	0.05	4.76	0.04	4.96	0.04	2.46
4	0.05	2.52	0.05	4.92	0.04	5.00	0.05	2.53
5	0.05	2.51	0.05	5.04	0.05	5.02	0.05	2.49
6	0.05	2.50	0.05	4.99	0.05	5.03	0.05	2.51
7	0.05	2.51	0.05	5.01	0.05	4.99	0.05	2.50
8	0.05	2.53	0.05	4.99	0.05	4.99	0.05	2.48
9	0.05	2.51	0.05	5.00	0.05	5.05	0.05	2.50
10	0.05	2.56	0.05	4.96	0.05	4.97	0.05	2.52
11	0.05	2.50	0.05	5.01	0.05	5.00	0.05	2.50

TABLE 11.2. Average SRC class residuals  $\text{err}_1(\mathbf{y}) := \|\mathbf{y} - X_{\text{tr}}\delta_1(\boldsymbol{\alpha}_{1,\epsilon})\|_2$  and  $\min_{2 \leq l \leq L}\{\text{err}_l(\mathbf{y}) := \|\mathbf{y} - X_{\text{tr}}\delta_l(\boldsymbol{\alpha}_{1,\epsilon})\|_2\}$  (over 1000 trials) on the random database model in the case of noise.

is smaller in this case, due to the class 1 training samples being uniformly distributed on  $S^{m-1}$ .

### 11.4. Summary

In this project, we designed a model, inspired by the work of Wright and Ma [102], for facial recognition and other similar classification databases. To model the mechanisms of SRC [104], we randomly generated a test sample as a non-negative linear combination of a single class's training samples. We computed the corresponding (sparse) coefficient vector and then ran experiments to test whether or not  $\ell^1$ -minimization, as it is used in the SRC setting, could recover this vector under increasing values of correlation, both within-class and in the database as a whole.

The results demonstrate that the within-class correlation in this model consistently improves  $\ell^1/\ell^0$ -recovery when compared to randomly-generated uniform data on the sphere. This is an important empirical result, as this latter type of data is one of the “golden children” of  $\ell^1/\ell^0$ -equivalence; i.e., these type of dictionaries produce, in some sense, ideal recovery (see, e.g., the work of Donoho [29]). However, those results are strongly asymptotic, and our experiments dealt only with small databases. More work is needed to determine if our findings hold up on larger data sets.

It is not too surprising, given the mutual coherence recovery condition studied in the last chapter, that very large correlation in the database as a whole can degrade recovery. When the

## 11.4. SUMMARY

---

global correlation in our model was very high, so that the classes, or sub-bouquets, began to overlap, we saw that  $\ell^1$ -minimization did not find the correct support of the sparse solution. However, we showed that the support could be completely fixed by a simple thresholding technique.

We also demonstrated that  $\ell^1$ -minimization achieved a good approximation of the sparsest solution in the case of noise in our model. Though the accuracy of the approximation generally decreased as the data became more correlated, this deterioration was slow compared to the increase in mutual coherence of the database. Further, the amount of  $\ell^2$ -error appeared to be less dependent on the relationship between noise  $\zeta$  and error tolerance  $\epsilon$  than it was on the amount of redundancy in the database.

Assuming that test samples truly are linear combinations of their ground truth class training samples, as is done in SRC, these experiments suggest that  $\ell^1$ -minimization will recover this class representation, leading to good classification in SRC and similar classification algorithms. This of course assumes that our model is appropriate for the given data set, and that its values of  $N_0$ ,  $m$ , and  $L$  are comparable to those used in our experiments, so that the class representation is sparse.

Our results are purely empirical; however, they strongly suggest that theoretical recovery results are possible. We conjecture that exact recovery can be provably obtained whenever the classes are sufficiently non-overlapping and that a similar result can be obtained in the case of noise. The amount of redundancy in the database and the number of classes will play a crucial role in this analysis.

Finally, though this project explicitly modeled the cone structure of facial images, our results are likely applicable to other areas of classification as well. In particular, as long as it is assumed that the training samples within each class are highly-correlated, we could amend our model so that the sign of each training sample was chosen randomly and so that the test sample was generated in the linear (not necessarily positive) span of its same-class training samples. However, since  $\ell^1$ -minimization is invariant to multiplication of the dictionary elements by  $\pm 1$ , we largely suspect that our results would be the same.

## CHAPTER 12

# Proving Equivalence in SRC via Nonlinear Embedding

### 12.1. This Approach

In this chapter, we consider *forcing* the mutual coherence criterion of  $\ell^1/\ell^0$ -equivalence to hold via data transformation. We want to analyze the effect of obtaining the *provably* sparsest solution in SRC on the algorithm’s classification accuracy.

As we have seen in the previous two chapters, there is an innate conflict in using  $\ell^1$ -minimization to obtain the sparsest solution in this context: we want the classes to be well-clustered, so that they may be easily classified, but if the data samples are not well-spread out, then we cannot prove that  $\ell^1/\ell^0$ -equivalence holds. As we have seen, class structure often results in large mutual coherence of the data set, making it impossible to apply either of the mutual coherence equivalence guarantees given in Theorem 9.2.4 and Theorem 9.2.5. We consider a resolution to this conflict through the use of *more space*. That is, if we had many “extra” dimensions, the data in each class could conceivably be spread out and we would still have enough “room” to keep the classes well-separated from each other.

Let us illustrate this in low dimension. Consider the toy example in which we have  $L = 2$  classes, each containing  $N_0 = 2$  samples in  $\mathbb{R}^m$  for  $m = 2$ . First, let the goal be to arrange the samples in a way that minimizes their mutual coherence while at the same time provides some indication of class. Assuming that the samples must be normalized (as in SRC), this class-structure criterion can reasonably be interpreted as the requirement that

$$\left| \langle \mathbf{x}_i^{(1)}, \mathbf{x}_j^{(1)} \rangle \right| > \left| \langle \mathbf{x}_i^{(1)}, \mathbf{x}_j^{(2)} \rangle \right| \text{ and } \left| \langle \mathbf{x}_i^{(2)}, \mathbf{x}_j^{(2)} \rangle \right| > \left| \langle \mathbf{x}_i^{(2)}, \mathbf{x}_j^{(1)} \rangle \right|,$$

for  $i, j \in \{1, 2\}$ . In other words, the samples in the same class must be more correlated than samples in different classes.

## 12.1. THIS APPROACH

---

One solution is given by the class matrices

$$X^{(1)} = [\mathbf{x}_1^{(1)}, \mathbf{x}_2^{(1)}] = \begin{bmatrix} 1 & \cos(\frac{\pi}{4} - \epsilon) \\ 0 & \sin(\frac{\pi}{4} - \epsilon) \end{bmatrix},$$

$$X^{(2)} = [\mathbf{x}_1^{(2)}, \mathbf{x}_2^{(2)}] = \begin{bmatrix} 0 & \cos(\frac{3\pi}{4} - \epsilon) \\ 1 & \sin(\frac{3\pi}{4} - \epsilon) \end{bmatrix}.$$

The (magnitude of) the dot product between samples in the same class is  $\cos(\frac{\pi}{4} - \epsilon)$ , and that of samples in different classes is  $\cos(\frac{\pi}{4} + \epsilon)$ . Clearly, the former quantity is the mutual coherence of the data set. This arrangement is illustrated in Figure 12.1a with  $\epsilon = 0.2$ .

Now, consider the same problem but in the case that we are given a third dimension. It is clear that we will be able to decrease the mutual coherence of the data set by moving samples into this extra space. One solution is given by the class matrices

$$X^{(1)} = \begin{bmatrix} 1 & \cos(\theta_1) \sin(\phi_1) \\ 0 & \sin(\theta_1) \sin(\phi_1) \\ 0 & \cos(\phi_1) \end{bmatrix},$$

$$X^{(2)} = \begin{bmatrix} 0 & \cos(\theta_2) \sin(\phi_2) \\ 1 & \sin(\theta_2) \sin(\phi_2) \\ 0 & \cos(\phi_2) \end{bmatrix},$$

for  $\theta_1 = \pi/4 - \epsilon$ ,  $\theta_2 = \pi/4 + \epsilon$ ,  $\phi_1 = 3\pi/4$ , and  $\phi_2 = \pi/4$ . The mutual coherence of the data set is  $\cos(\frac{\pi}{4} - \epsilon) \sin(\frac{3\pi}{4}) = \sin(\frac{\pi}{4} + \epsilon) \cos(\frac{\pi}{4})$ . This arrangement is illustrated in Figure 12.1b with  $\epsilon = 0.2$ .

For  $\epsilon = 0.2$ , for example, adding an additional dimension allows us to decrease the mutual coherence of the data set from  $\cos(\frac{\pi}{4} - \epsilon) \approx 0.8335$  to  $\cos(\frac{\pi}{4} - \epsilon) \sin(\frac{3\pi}{4}) \approx 0.5894$ . This is a substantial decrease.

The reader may question the practical relevance of considering such a transformation. We admit that our motivation behind this project is mostly theoretical: to what extent does finding the sparsest solution in SRC affect classification? Is the reward of true sparsity worthwhile, or is it enough in the SRC framework to achieve merely approximately-sparse solutions? Without knowing

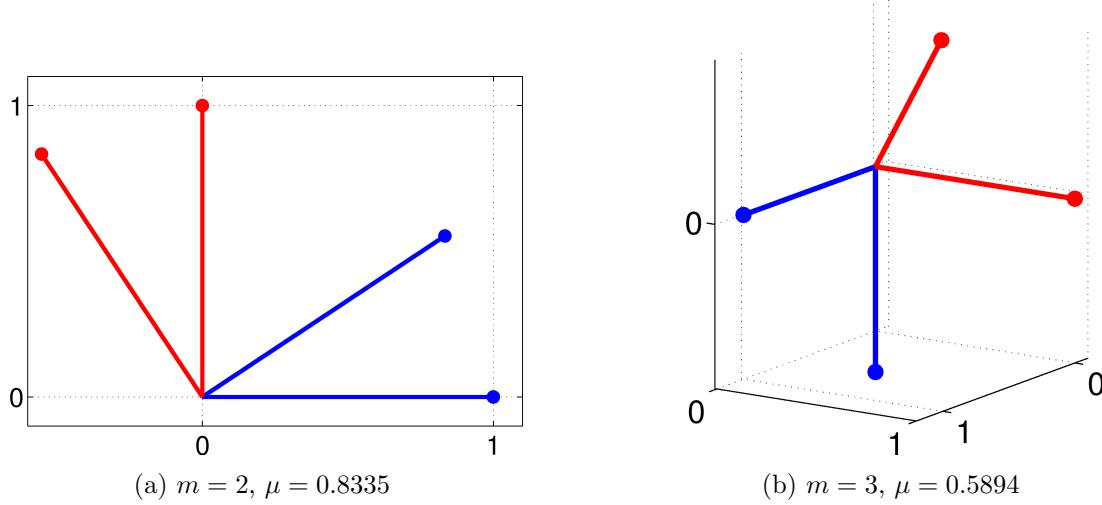


FIGURE 12.1. Illustration of decreasing the mutual coherence of a data set by embedding it into higher dimension. (a) Samples in original space  $\mathbb{R}^2$ , (b) Samples in transform space  $\mathbb{R}^3$ . Colors denote classes.

the sparsest solution, these questions cannot be answered. Rather than using a combinatorial solver to determine the sparsest solution  $\alpha_0$  (an impractical undertaking for large data sets), we consider a transformation of the data that allows us to apply Theorem 9.2.4.

## 12.2. Designing the Transform

**12.2.1. Considering an Explicit Transform.** Ideally, we seek a transform  $\phi^* : \mathbb{R}^m \rightarrow \mathbb{R}^{\tilde{m}}$  with  $m < \tilde{m}$  that satisfies

$$\phi^* = \arg \max_{\phi \in \mathcal{C}} f_{\text{cs}}([\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_{N_{\text{tr}}})]) \text{ subject to } \mu([\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_{N_{\text{tr}}})]) \leq \tilde{\mu},$$

where  $\mathcal{C}$  is some compact set (so that  $f_{\text{cs}}$  obtains a maximum). Here,  $f_{\text{cs}}$  evaluates the amount of class structure in the transformed database (“cs” stands for “class structure”), and  $\tilde{\mu}$  is an upper bound on the transformed database’s mutual coherence. For example,  $f_{\text{cs}}$  might denote the inverse of the sum of within-class distances or the inverse of the Frobenius norm of the within-class scatter matrix used in *linear discriminant analysis* [42, 74].

We note that  $\phi^*$  must be a nonlinear transform, otherwise the dimension of the subspace containing the embedded samples will be no greater than that of the original space ( $m$ ). Thus we

## 12.2. DESIGNING THE TRANSFORM

---

will have failed to utilize the extra space (needed to achieve our objective) awarded by the increased ambient dimension  $\tilde{m}$ .

**12.2.2. Using the Kernel Method.** Rather than constructing an explicit transform, we consider the reduction of mutual coherence via the so-called “kernel trick.” Essentially, the kernel trick allows us to perform operations in a space of dimension  $\tilde{m} > m$  (possibly infinite-dimensional) without having to actually compute the transformed samples. The “trick” is to work only with the inner-products between transformed samples, which are given to us by some kernel function  $\kappa : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ .

More formally, denote the transform by  $\phi$ . We will write  $\Phi(X_{\text{tr}}) := [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_{N_{\text{tr}}})]$  to denote the matrix of training samples in the transform (or *kernel*) space. We define the inner-product in the kernel space as

$$\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle =: \kappa(\mathbf{x}_i, \mathbf{x}_j),$$

for  $1 \leq i, j \leq N_{\text{tr}}$ . The kernel function  $\kappa$  should satisfy *Mercer’s condition*<sup>1</sup> so that  $\kappa$  defines a proper inner-product [22].

Kernel methods can be particularly effective when used to “non-linearize” linear classifiers. In *kernel support vector machines*, for example, classes that are not linearly-separable in the original space may be separated linearly in kernel space (see the work of Boser et al. [7]). Though SRC is not linear, it does assume a linear relationship between the test sample and the training samples in its ground truth class. When such a relationship does not hold in the original space, it may hold in kernel space given that an appropriate kernel is selected [111].

The kernel trick may seem like magic, but it has a drawback. Namely, the linear structure that we desire in kernel space cannot be attained by just any kernel. Kernel methods require some prior knowledge about the data and the classification task in order to select a kernel that will lead to the desired structure in kernel space.

**12.2.3. The Gaussian Kernel.** We consider the *Gaussian* kernel, given by

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) := e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}}.$$

---

<sup>1</sup>The kernel  $\kappa$  satisfies Mercer’s condition if  $\iint \kappa(\mathbf{x}, \mathbf{y}) g(\mathbf{x})g(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0$  for all square-integrable functions  $g$ .

## 12.2. DESIGNING THE TRANSFORM

---

Essentially, the Gaussian kernel adds *inverse exponential scaling* to the Euclidean distance function. Points close together obtain values of  $\kappa$  that are close to 1, whereas points that are faraway from each other have kernel values approaching 0. The *window* or *width* parameter  $\sigma$  controls the drop off (or steepness) of this trade-off.

Recall that our goal in the transform space is to decrease mutual coherence while still maintaining some degree of class structure in the data. In kernel space, the mutual coherence of the (transformed) training set is given by

$$\begin{aligned}\mu(\Phi(X_{\text{tr}})) &= \max_{1 \leq i \neq j \leq N_{\text{tr}}} |\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle| \\ &= \max_{1 \leq i \neq j \leq N_{\text{tr}}} |\kappa(\mathbf{x}_i, \mathbf{x}_j)| \\ &= \max_{1 \leq i \neq j \leq N_{\text{tr}}} e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}}.\end{aligned}$$

Since the vectors  $\hat{\mathbf{x}}_i$  and  $\hat{\mathbf{x}}_j$  satisfying  $\|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j\|_2 = \max_{i \neq j} \|\mathbf{x}_i - \mathbf{x}_j\|_2$  are fixed for a given training set, the mutual coherence  $\mu(\Phi(X_{\text{tr}}))$  depends completely on  $\sigma$ . Thus we can write  $\mu(\Phi(X_{\text{tr}})) = \mu = \mu(\sigma)$ . To reiterate, we can completely control the mutual coherence of the data in the kernel space by adjusting  $\sigma$ . Thus the Gaussian kernel appears to be a good choice for use in our task, assuming that we also have some means of controlling the amount of class structure in the data set in kernel space.

How do we help ensure that class structure is maintained as the mutual coherence decreases? As  $\sigma \rightarrow 0$ ,  $\mu(\sigma) \rightarrow 0$ . In kernel space, this means that the training samples become closer to being pairwise orthogonal, the result of the inner-products between samples in kernel space—equivalently, their Gaussian kernel values in the original space—going to 0. Although the classes eventually become indistinguishable as  $\sigma \rightarrow 0$ , samples that are close together (with respect to Euclidean distance) in the original space take longer to become pairwise orthogonal than samples that are far apart. It follows that data sets that will maintain their class structure better for small values of  $\sigma$  are exactly those that are well-separated (in the original space) by Euclidean distance. Thus the class structure criterion can be satisfied by choosing a database with this property.

This setup might raise some objections. In the first place, were our aim simply to classify a given data set, this line of approach would be circular. If we knew the data were easily classified

## 12.2. DESIGNING THE TRANSFORM

---

using Euclidean methods, we could simply apply an algorithm such as 1NN in the original space. Secondly, often data sets cannot be well-separated using Euclidean distance, especially when the data from each class is drawn from a nonlinear manifold.

However, our aim in this experiment is *not* to classify an arbitrary data set or to validate the efficacy of the kernel method for use in classification. As we expand upon in the next section, our goal is to *leverage the kernel setup* in order to investigate the relationship between the mutual coherence  $\ell^1/\ell^0$ -equivalence guarantee, the sparsity of the coefficient vector, and the classification accuracy of SRC. It is a tool to access the *provably sparsest solution*, something that we cannot (tractably) do in general in the original space. Further, remember that the key to successfully using kernel methods is to select a kernel designed to produce the desired structure in kernel space. For our purposes, this is the *preservation of class structure as the samples in kernel space become less correlated*. When the Gaussian kernel is used, this exactly means that Euclidean distance is a good indicator of class in the original space.

**12.2.4. The Goal.** Our goal is to use the kernel trick with the Gaussian kernel to achieve  $\ell^1/\ell^0$ -equivalence in kernel space. We want to investigate the effect of this (provable) equivalence on the classification accuracy of SRC.

In order to ensure  $\ell^1/\ell^0$ -equivalence, we will choose the Gaussian width parameter  $\sigma$  so that the mutual coherence is small enough that Theorem 9.2.4 holds. Let us set

$$k_{\sup} := \frac{1}{2} \left( 1 + \frac{1}{\mu} \right).$$

Clearly,  $k_{\sup}$  completely depends on  $\mu$ , or equivalently, on  $\sigma$ . As  $\sigma$  approaches 0,  $k_{\sup} = k_{\sup}(\sigma)$  blows up. Since class structure deteriorates as  $\sigma \rightarrow 0$  (as samples become more and more pairwise orthogonal), we want to choose  $\sigma$  to be the largest value such that, with high probability (whp), the sparsity level  $\|\boldsymbol{\alpha}_1\|_0$  is less than  $k_{\sup}$ , where  $\boldsymbol{\alpha}_1 = \boldsymbol{\alpha}^* \in \mathbb{R}^{N_{\text{tr}}}$  is the solution to the exact  $\ell^1$ -minimization problem in SRC given by Eq. (3.5) (replacing  $X_{\text{tr}}$  with  $\Phi(X_{\text{tr}})$  and  $\mathbf{y}$  with  $\phi(\mathbf{y})$ ). This will ensure that  $\boldsymbol{\alpha}_1$  is the sparsest solution by Theorem 9.2.4. Using “mc” to denote “mutual coherence,” we define

$$(12.1) \quad \sigma_{\text{mc}} := \max \left\{ \sigma : \|\boldsymbol{\alpha}_1\|_0 \stackrel{\text{whp}}{<} k_{\sup} \right\}.$$

### 12.3. AN UNFORTUNATE YET NECESSARY MODIFICATION

---

Consider now what happens as  $\sigma \rightarrow \infty$ . We have that  $\mu \rightarrow 1$  and  $k_{\text{sup}}(\sigma) \rightarrow 1$ . More informatively, the samples in the kernel space become increasingly correlated, and class structure is eventually lost. Another way of putting this is that as  $\sigma \rightarrow 0$ , “within-class” structure is lost, and as  $\sigma \rightarrow \infty$ , “between-class” structure is lost. In both cases, the data becomes harder to classify using SRC-based classifiers. It follows, as for all kernel classifiers when the Gaussian kernel is used, that there is a set of values of  $\sigma$  (possibly containing only one point) that produces a sufficiently high degree of class structure so that maximum classification accuracy is achieved for all values in this set (whp). We denote the maximum value in this set by  $\sigma_{\text{acc}}$ .

We want to investigate the relationship between  $\sigma_{\text{mc}}$  and  $\sigma_{\text{acc}}$ . We are also interested in the sparsity level  $\|\boldsymbol{\alpha}_1\|_0$  at both  $\sigma = \sigma_{\text{mc}}$  and  $\sigma = \sigma_{\text{acc}}$ . Since some coefficients may be small, we lastly should consider the size of the coefficients of training samples corresponding to the ground truth class of  $\mathbf{y}$ . In analyzing these quantities and relationships, we aim to provide insight into the role of sparsity in classification.

### 12.3. An Unfortunate yet Necessary Modification

On a toy database, we confirmed that the kernel method can work exactly as we want it to, producing training samples in kernel space that have small enough within-class correlation so that the bound in Eq. (9.4) holds, while still distinguishing classes via even smaller between-class coherence. For very well-separated data (in the original space), we found that when  $\sigma = \sigma_{\text{mc}}$ , the coherence between training samples in different classes, and that between the test sample and training samples not in its ground truth class, were actually zero (in kernel space). It seems that this should allow for very good classification by SRC in kernel space *and* provable  $\ell^1/\ell^0$ -equivalence of the solution vector, allowing us to investigate the change in classification accuracy as  $\sigma$  increases and the solution becomes denser.

Unfortunately, as  $\sigma \rightarrow 0$ , we lose the ability to write  $\Phi(X_{\text{tr}})\boldsymbol{\alpha} = \phi(\mathbf{y})$  for *any* coefficient vector  $\boldsymbol{\alpha}$ . Recall that this equality is a key aspect of the mutual coherence condition in Theorem 9.2.4. In decreasing  $\sigma$ , we cause not only the training samples to become more orthogonal to each other, but also the test sample to become more orthogonal to each training sample, to the point that when

### 12.3. AN UNFORTUNATE YET NECESSARY MODIFICATION

---

$\sigma = \sigma_{\text{mc}}$ ,  $\phi(\mathbf{y})$  is likely not contained in the span of the columns of  $\Phi(X_{\text{tr}})$ . We state this formally in the following proposition:

PROPOSITION 12.3.1. *Given the kernel setup outlined above, suppose that we take  $\sigma$  small enough so that the mutual coherence bound  $\|\boldsymbol{\alpha}\|_0 < \frac{1}{2} \left(1 + \frac{1}{\mu(\Phi(X_{\text{tr}}))}\right)$  holds, i.e.,  $\sigma \leq \sigma_{\text{mc}}$ . Then*

$$\Phi(X_{\text{tr}})\boldsymbol{\alpha} \neq \phi(\mathbf{y})$$

*with high probability, for all  $\boldsymbol{\alpha} \in \mathbb{R}^{N_{\text{tr}}}$ .*

**Proof:** This is a direct consequence of Corollary 10.4.3.  $\square$

To summarize, a consequence of our kernel space setup is the problem of “too much space:” the resulting system is overdetermined with no solution to  $\Phi(X_{\text{tr}})\boldsymbol{\alpha} = \phi(\mathbf{y})$  when  $\sigma \leq \sigma_{\text{mc}}$ . By Theorem 9.2.4, the minimal  $\ell^1$ -norm solution satisfying  $\Phi(X_{\text{tr}})\boldsymbol{\alpha}_1 = \phi(\mathbf{y})$  with  $\|\boldsymbol{\alpha}_1\|_0 < (1/2)(1 + (1/\mu))$  is necessarily the sparsest such solution. However, if there is no solution satisfying  $\Phi(X_{\text{tr}})\boldsymbol{\alpha} = \phi(\mathbf{y})$ , then there can be no sparsest solution!

Even when the equality in SRC is relaxed and the constrained  $\ell^1$ -minimization problem in Eq. (9.3) is used,<sup>2</sup> relating the found solution  $\boldsymbol{\alpha}_{1,\epsilon}$  and the true sparsest solution  $\boldsymbol{\alpha}_0$  using Theorem 9.2.5 requires the *existence* of some  $\boldsymbol{\alpha} = \boldsymbol{\alpha}_0$  satisfying the equality  $\Phi(X_{\text{tr}})\boldsymbol{\alpha} = \phi(\mathbf{y})$ . Since the bound in Eq. (9.5) in the noisy case is more restrictive than Eq. (9.4) in the noiseless case, to satisfy Theorem 9.2.5 we must have  $\sigma < \sigma_{\text{mc}}$ . By Proposition 12.3.1, no such  $\boldsymbol{\alpha}$  exists, and it follows that Theorem 9.2.5 cannot be applied in this setup, either.

We can avoid this debilitating problem by generating  $\phi(\mathbf{y})$  as a linear combination of the columns of  $\Phi(X_{\text{tr}})$  (in particular, with nonzero coefficients occurring at training samples in the ground truth class of  $\mathbf{y}$ ). In this manner, we can ensure that  $\Phi(X_{\text{tr}})\boldsymbol{\alpha} = \phi(\mathbf{y})$  has a solution for all values of  $\sigma$ . However, this will mean that we never actually compute or handle the test sample  $\mathbf{y}$  in the original space and only *assume that it exists implicitly*, working solely with the inner-products between samples in kernel space. This is a little weird! The reader may object that we cannot just make up test samples in this manner, and in general, this is absolutely true. Nevertheless, we

---

<sup>2</sup>Note that the formulation in Eq. (9.3) is equivalent to the regularized  $\ell^1$ -minimization problem in Eq. (3.7) in the formal SRC algorithm statement.

### 12.3. AN UNFORTUNATE YET NECESSARY MODIFICATION

---

reiterate that our goal in this experiment is not to validate the kernel approach or to classify an arbitrary database, and so the implied existence of  $\mathbf{y}$  is acceptable in this context.

We admit, though, that this modification is not ideal, and it changes the experiment significantly. In fact, our whole motivation must change, as we no longer need to use Theorem 9.2.4 to find a sparse solution. For all values of  $\sigma$ , we now know that there exists a solution  $\boldsymbol{\alpha}$  to  $\Phi(X_{\text{tr}})\boldsymbol{\alpha} = \phi(\mathbf{y})$  having no more than  $\|\boldsymbol{\alpha}\|_0 = N_l$  nonzero coefficients, where  $\mathbf{y}$  is in class  $l$ .

Additionally, as  $\sigma \rightarrow 0$  and the training data become closer to orthogonal, we will never lose the relationship  $\phi(\mathbf{y}) \in \text{span}\{\phi(\mathbf{x}_1^{(l)}), \dots, \phi(\mathbf{x}_{N_l}^{(l)})\}$ . Thus we will not see the classification performance deteriorate *at all* as  $\sigma \rightarrow 0$ . In other words, decreasing  $\mu$  so that we can provably obtain  $\ell^1/\ell^0$ -equivalence in this setup can only *help* classification accuracy, as doing so isolates the linear relationship between the test sample and the training samples in its ground truth class (in kernel space). As a consequence, we do not need to explicitly require that the database be well-separated by Euclidean distance in the original space, since this imposed linear relationship will essentially do this for us. We stress that this is certainly not the case in general: consider the increasing difficulty of identifying class structure in a data set whose samples become more and more uncorrelated. Even if the classes are very well-separated in the original space, there will be a certain point at which all samples become orthogonal and class structure is entirely lost. Thus generating  $\phi(\mathbf{y})$  in this manner adds an undesirable degree of artificiality into our experiment.

However, we will show that we can still gain informative insights from observing the relationship between  $\sigma_{\text{mc}}$  and  $\sigma_{\text{acc}}$  in this setup. Further, we will observe the sparsity level of the coefficient vectors at  $\sigma_{\text{acc}}$  and investigate how  $\sigma_{\text{acc}}$  depends on the difficulty of the classification task. Note that these goals make this experiment fundamentally different from those in Chapter 11, wherein we primarily focused on the recovery of the sparsest solution using  $\ell^1$ -minimization and never actually performed classification using SRC. We also did not consider the relationship between classification accuracy and the mutual coherence bound in Eq. (9.4) in that project. Thus our work here will provide new insights.

## 12.4. Experiments

**12.4.1. Experimental Setup.** For a fixed training set and fixed  $\sigma$ , we generate  $N_l$  test samples in kernel space for each class  $1 \leq l \leq L$  as linear combinations of the training samples in that class (in kernel space) with coefficients randomly drawn from  $\text{unif}(0, 1)$  distribution. Non-negative coefficients are used so that  $\langle \phi(\mathbf{y}), \phi(\mathbf{x}_j) \rangle \geq 0$  for  $1 \leq j \leq N_{\text{tr}}$ , as is consistent with the Gaussian kernel. We then apply SRC in kernel space to classify the resulting test samples, using the Kernel SRC algorithm of Kang et al., in particular, their *kernel coordinate descent* (KCD) algorithm [55]. Note that in their paper, the authors apply this algorithm to the *local binary patterns* of the original samples instead of the original samples themselves, and since other types of kernels are more appropriate for these type of features, they do not use the Gaussian kernel, as we do.

In our experiments, we determine  $\sigma_{\text{mc}}$  and  $\sigma_{\text{acc}}$  by trial-and-error. Given the randomness inherent in the database construction (see Section 12.4.3 for a description of the database used), we request that the reader grant us leeway to make judicious choices in terms of rounding, etc., in determining these values. Additionally, note that we thresholded the entries of each  $\ell^1$ -minimized coefficient vector  $\boldsymbol{\alpha}_1$  by  $10^{-10}$  to standardize machine precision.

**12.4.2. An Upper Bound.** We saw in Chapter 10 that we cannot apply Theorem 9.2.4 unless  $\mu < \frac{1}{3}$ . Since we are using the kernel approach, this means that we must have

$$\mu(\Phi(X_{\text{tr}})) = \max_{1 \leq i \neq j \leq N_{\text{tr}}} \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \max_{1 \leq i \neq j \leq N_{\text{tr}}} \kappa(\mathbf{x}_i, \mathbf{x}_j) < \frac{1}{3}.$$

In particular, since we are using the Gaussian kernel, it must be the case that

$$\begin{aligned} \max_{i \neq j} \kappa(\mathbf{x}_i, \mathbf{x}_j) &= \max_{1 \leq i \neq j \leq N_{\text{tr}}} e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}} < \frac{1}{3} \\ \Rightarrow \sigma &< \frac{1}{\sqrt{\log(3)}} \max_{1 \leq i \neq j \leq N_{\text{tr}}} \|\mathbf{x}_i - \mathbf{x}_j\|_2. \end{aligned}$$

Since the training samples (in the original space) are normalized, this means that

$$(12.2) \quad \sigma < \frac{2}{\sqrt{\log(3)}} \approx 1.35.$$

Thus in searching for  $\sigma_{\text{mc}}$ , we only need to consider values of  $\sigma$  less than 1.35.

## 12.4. EXPERIMENTS

---

**12.4.3. Database Description.** We constructed a very simple toy database in the original space as follows: Samples in the  $l$ th class were initially  $N_0$  copies of the canonical basis vector  $\mathbf{e}_l \in \mathbb{R}^L$ , where  $L$  was the number of classes. The feature dimension  $m$  was user-specified, and then  $m - L$  coordinates were added to each canonical basis vector and set to zero. Lastly, random noise from  $\mathcal{N}(0, \eta^2)$  was added to all (training) samples in all coordinates.

We set  $N_0 = 5$ ,  $m = 50$ , and  $L = 20$ , so that each class would consist of a relatively small portion of the dictionary  $X_{\text{tr}}$ , as is ideal in SRC. We used three different values of noise level  $\eta \in \{0.001, 0.1, 0.5\}$ . As in the  $\ell^1$ -minimization algorithm HOMOTOPY, KCD requires an error/sparsity tradeoff parameter  $\lambda$ . To force near-exactness in the representations, we set  $\lambda = 10^{-10}$ .

Note that one advantage of the kernel setup in this context is the ability to apply it to *any* training database. This is another critical difference between this project and the experiments in Chapter 11, which were fundamentally based around a particular database model.

**REMARK 12.4.1.** *The reader may question why, in this project, we did not use the synthetic database from the experiments in Chapters 5 and 7, or why we used a different synthetic database than the one in the last chapter, so that we might obtain a fair comparison. We are making our best effort to stress that this line of thinking misconstrues the point of this experiment. Here, we only care about the classification results of Kernel SRC as they relate to the sparsity level  $\|\boldsymbol{\alpha}_1\|_0$  and the mutual coherence bound in Eq. (9.4). We are completely uninterested in whether the kernel approach improves the classification accuracy of SRC (for a positive answer to this question, see, for example, Kang et al.’s paper [55]). Further, the aforementioned synthetic databases were designed for specific purposes: the sinusoidal wave database used in the LPCA-SRC experiments was constructed so that LPCA-SRC and SRC could be compared in the case of nonlinear and intersecting class manifolds, and the database model in Chapter 11 was chosen so that  $\ell^1/\ell^0$ -equivalence—and not, specifically, classification performance in SRC—could be studied at increasing levels of data correlation. So that the aims of these previous experiments did not bleed into our goals here, we used a completely new (and very simple) database.*

**12.4.4. Results.** In Figure 12.2, we plot the synthetic database results for each value of  $\eta$  over various values of  $\sigma$ , annotating the values of  $\sigma_{\text{mc}}$  and  $\sigma_{\text{acc}}$ . We report the averages over 100

## 12.4. EXPERIMENTS

---

instantiations of the training and test sets (“trials”). In particular, we report the average sparsity level, Kernel SRC classification accuracy, and the (relative)  $\ell^2$  and  $\ell^1$ -norms of the correct class support. These quantities are defined rigorously as

$$\text{Sparsity} := \underset{\text{all trials}}{\text{mean}} \left\{ \underset{\text{all test samples}}{\text{median}} \frac{\|\boldsymbol{\alpha}_1\|_0}{N_{\text{tr}}} \right\}$$

for  $\boldsymbol{\alpha}_1$  thresholded at  $10^{-10}$  (we compute the median sparsity over all test samples so that the result is more robust to atypical very sparse or very dense coefficient vectors),

$$\text{Accuracy} := \underset{\text{all trials}}{\text{mean}} \left\{ \underset{\text{all test samples}}{\text{mean}} \mathbb{1}_{\{\text{class\_label}(\mathbf{y}) = \text{ground\_truth\_class}(\mathbf{y})\}} \right\}$$

where  $\mathbb{1}_{\{x=y\}}$  is the indicator function that returns 1 if  $x = y$  and 0 otherwise, and

$$\text{supp}(\ell^2) := \underset{\text{all trials}}{\text{mean}} \left\{ \underset{\text{all test samples}}{\text{mean}} \frac{\|\delta_{\text{GT}}(\boldsymbol{\alpha}_1)\|_2}{\|\boldsymbol{\alpha}_1\|_2} \right\}, \quad \text{supp}(\ell^1) := \underset{\text{all trials}}{\text{mean}} \left\{ \underset{\text{all test samples}}{\text{mean}} \frac{\|\delta_{\text{GT}}(\boldsymbol{\alpha}_1)\|_1}{\|\boldsymbol{\alpha}_1\|_1} \right\},$$

where the nonzero entries of  $\delta_{\text{GT}}(\boldsymbol{\alpha}_1)$  are exactly those from  $\boldsymbol{\alpha}_1$  that correspond to the ground truth class of the given test sample.

From Figure 12.2, we see that  $\sigma_{\text{acc}}$  was generally much larger than  $\sigma_{\text{mc}}$ , and that the Kernel SRC method could tolerate substantial  $\ell^1$  and  $\ell^2$ -support error before classification deteriorated. Further, perfect classification was achieved even for maximally dense  $\boldsymbol{\alpha}_1$ . This suggests that a strictly-sparse solution vector is not necessary to the success of SRC.

As the level of noise  $\eta$  increased, we see in Figure 12.2 that  $\sigma_{\text{acc}}$  decreased towards  $\sigma_{\text{mc}}$ . However,

$$\lim_{\eta \rightarrow \infty} \sigma_{\text{acc}} \neq \sigma_{\text{mc}}.$$

Once the class structure was lost due to noise in the original space, increasing the noise level further had no effect on the quantities displayed in Figure 12.2. In other words, Figure 12.2c is representative of the results for larger values of  $\eta$ .

We also observe that for  $\eta = 0.001$ , the sparsest solution was obtained by  $\ell^1$ -minimization for values of  $\sigma$  slightly larger than  $\sigma_{\text{mc}}$  (note the position of the  $\sigma_{\text{mc}}$  arrow tip in Figure 12.2a). In fact, the mutual coherence of the data set with  $\eta = 0.001$  reached  $\mu = 0.9994$  before  $\ell^1$ -minimization failed to retrieve the sparsest solution. This indicates that when the classes are well-separated (for small  $\eta$  and sufficiently small  $\sigma$ , separability in the original space carries over to kernel space in

## 12.4. EXPERIMENTS

---

this experiment),  $\ell^1/\ell^0$ -equivalence can still be achieved even when the mutual coherence is much larger than that allowed by Eq. (9.4). This reinforces the findings from the last chapter, namely, that  $\ell^1/\ell^0$ -equivalence holds on highly-correlated data as long as the vectors corresponding to the support of the sparsest solution are sufficiently separated from the other dictionary elements. On the other hand, for larger values of  $\eta$ , i.e., when the classes were less well-separated, the bound in Eq. (9.4) appears to be approximately tight.

**12.4.5. Examining the Accuracy Threshold.** It is notable that  $\sigma_{\text{acc}}$  is substantially larger than  $\sigma_{\text{mc}}$  for all  $\eta$ , and that the accuracy in Kernel SRC has a steep drop-off as soon as  $\sigma > \sigma_{\text{acc}}$ . The value  $\sigma_{\text{acc}}$  appears to be a threshold for which the linear relationship between  $\phi(\mathbf{y})$  and the training samples in its ground truth class cannot be identified by the classification mechanism in (Kernel) SRC. We want to know what triggers this threshold.

We first look for an “elbow” or sharp change in the correlation between  $\phi(\mathbf{y})$  and training samples in its ground truth class, and that between  $\phi(\mathbf{y})$  and samples in other classes. In particular, we computed

$$\text{corr}_{\text{GT}} := \text{mean} \left\{ \underset{\text{all trials}}{\text{median}} \left\{ \underset{\mathbf{x}_j^{(l)}: \mathbf{y} \in \text{class } l}{\text{median}} \left\langle \phi(\mathbf{y}), \phi(\mathbf{x}_j^{(l)}) \right\rangle \right\} \right\}$$

and

$$\text{corr}_{\text{other}} := \text{mean} \left\{ \underset{\text{all trials}}{\text{median}} \left\{ \underset{l: \mathbf{y} \notin \text{class } l}{\text{median}} \left\{ \underset{1 \leq j \leq N_l}{\text{median}} \left\langle \phi(\mathbf{y}), \phi(\mathbf{x}_j^{(l)}) \right\rangle \right\} \right\} \right\}.$$

Again, we compute the median quantities within each trial to make the correlation values more robust to sample outliers.

The results for  $\eta = 0.1$  are shown in Figure 12.3. The plots for the other values of  $\eta$  are similar. As we can see, the accuracy threshold  $\sigma_{\text{acc}}$  occurred *after* the sharp increase in the correlation quantities. In fact, we see that SRC was able to retrieve the correct classification assignment when  $\text{corr}_{\text{GT}}$  was only moderately larger than  $\text{corr}_{\text{other}}$ . On the other hand, the sharp increase in the correlation quantities appears to correspond to the steep increase in sparsity level, which makes sense in the context of the mutual coherence equivalence guarantee in Theorem 9.2.4.

## 12.4. EXPERIMENTS

---

As a more informative approach to understanding the accuracy threshold, in particular, what causes the sharp drop-off in accuracy at  $\sigma_{\text{acc}}$ , we consider the distribution of the absolute values of the coefficients, i.e., the magnitude of the coordinates of  $\alpha_1$ , with respect to the different classes. In particular, for  $\eta = 0.1$ , we computed the mean vector  $|\alpha_1|$  over the  $N_0 = 5$  class  $l = 20$  training samples, and then averaged the result over 100 trials:

$$\underset{\text{all trials}}{\text{mean}} \left\{ \underset{\mathbf{y} \in \text{class } l=20}{\text{mean}} \{ |\alpha_1| \} \right\}.$$

We lastly normalized the resulting vector so that its entries summed to 1.

We plot the results in Figure 12.4 for a handful of representative values of  $\sigma$ . The  $x$ -axis in the left-hand-side plots (Figures 12.4a, 12.4c, 12.4e, and 12.4g) corresponds to the individual coordinates of the averaged vector  $|\alpha_1| \in \mathbb{R}^{N_{\text{tr}}}$ . The coordinates corresponding to training samples in each class are simply summed to produce the right-hand-side plots (Figures 12.4b, 12.4d, 12.4f, and 12.4h), so that the contribution from each class in the representation of  $\phi(\mathbf{y})$  can be viewed easily. We also include the corresponding Kernel SRC classification accuracies for reference.

Given the dominance of coefficients corresponding to class  $l = 20$  in Figures 12.4a-12.4d, it is not surprising that Kernel SRC obtains perfect accuracy in these cases. It is also quite clear from these figures that small coefficients in the wrong class do not negatively affect classification accuracy. Thus there is no reason to require a solution sparser than that with  $\sigma = 3$ .

For  $\sigma \in \{5, 9\}$ , the closeness in the coefficient magnitudes between those corresponding to class  $l = 20$  and those corresponding to other classes illustrates the decreased accuracy in Kernel SRC; recall that these plots contain averages. Additionally, we note that the distribution of the coefficients in class  $l = 20$  became fairly unbalanced among that class's training samples for these large values of  $\sigma$ . This is because as mutual coherence increased, the class  $l = 20$  samples became more and more parallel to each other. Thus most of  $\phi(\mathbf{y})$  could be represented using only the first training sample in that class.

Figure 12.4 helps to explain the sharp drop-off in accuracy at  $\sigma_{\text{acc}}$ . Though the quantities  $\text{corr}_{\text{GT}}$  and  $\text{corr}_{\text{other}}$  are only slightly increasing at  $\sigma_{\text{acc}}$  (and the general behavior of the coefficients varying smoothly), the threshold occurs right at the point that the coefficients of other classes become competitive with those from the correct class (as we would expect). The sharp drop-off can be

## 12.5. KEY FINDINGS

---

attributed to the nonlinearity of the `min` function in determining  $\min_{1 \leq l \leq L} \{\|\phi(\mathbf{y}) - \Phi(X_{\text{tr}})\delta_l(\boldsymbol{\alpha}_1)\|_2\}$  in the classification stage of (Kernel) SRC.

### 12.5. Key Findings

First of all, we acknowledge that this experiment did not satisfy all of the goals laid out in its introduction. However, we can still draw some important conclusions:

- Any procedure that spreads out the data in each class in a way that decreases mutual coherence yet aims to maintain class structure will necessarily come into conflict with maintaining a linear relationship between  $\mathbf{y}$  and *any* subset of training samples. More precisely, it is generally impossible to write  $\mathbf{y}$  as a linear combination of the training samples in class  $l$  while satisfying the bound

$$\|\boldsymbol{\alpha}\|_0 < \frac{1}{2} \left( 1 + \frac{1}{\mu(X_{\text{tr}})} \right) \approx \frac{1}{2} \left( 1 + \frac{1}{\mu([X^{(l)}, \mathbf{y}])} \right),$$

i.e., when  $\mathbf{y}$  is spread out in the same manner as the other samples in the database. Besides artificially generating  $\mathbf{y}$  as a linear combination of the training samples *after* they have been spread out, it is not clear to us how to overcome this conflict.

- Though generating  $\mathbf{y}$  as a linear combination of its ground truth class training samples in kernel space prevented us, in some sense, from isolating the relationship between  $\sigma_{\text{mc}}$  and classification accuracy, we were still able to study the correspondence between  $\sigma_{\text{mc}}$  and sparsity level  $\|\boldsymbol{\alpha}_1\|_0$ . In particular, we confirmed our previous findings that perfect recovery can be achieved on highly-correlated data as long as the classes are sufficiently well-separated (in this experiment, this meant small  $\eta$ ).
- We saw that there was a sharp drop-off in classification accuracy at  $\sigma_{\text{acc}} > \sigma_{\text{mc}}$ , which was not directly correlated with a sharp change in either sparsity or the relationship between within-class and between-class correlation, or in the normalized  $\ell^2$  and  $\ell^1$ -norms of  $\delta_{\text{GT}}(\boldsymbol{\alpha}_1)$ . Though  $\ell^1/\ell^0$ -equivalence (whether provable by Theorem 9.2.4 or not) was a way to ensure perfect classification accuracy in this experiment, it was not necessary. The classification mechanism in SRC can clearly tolerate even the maximal number of nonzero coefficients in the representation, as long as the magnitudes of coefficients corresponding

## 12.5. KEY FINDINGS

---

to the wrong classes are small with respect to those from the correct class. In this sense, *relative*—or *approximate*—sparsity is the key to SRC. It might be possible to make this idea precise in terms of a coefficient thresholding procedure similar to the one used in Chapter 11.

In future research, it would be interesting to consider the modification of the above experiment when noise is added to the test sample  $\phi(\mathbf{y})$  after it is generated as a linear combination of its ground truth class training samples in kernel space. Of course, this will not have the same effect as adding noise to the original (and implicitly-defined) test sample  $\mathbf{y}$ , but it would allow us to investigate the relationship between classification accuracy in SRC and the mutual coherence bound in the case of noise as stated in Theorem 9.2.5.

## 12.5. KEY FINDINGS

---

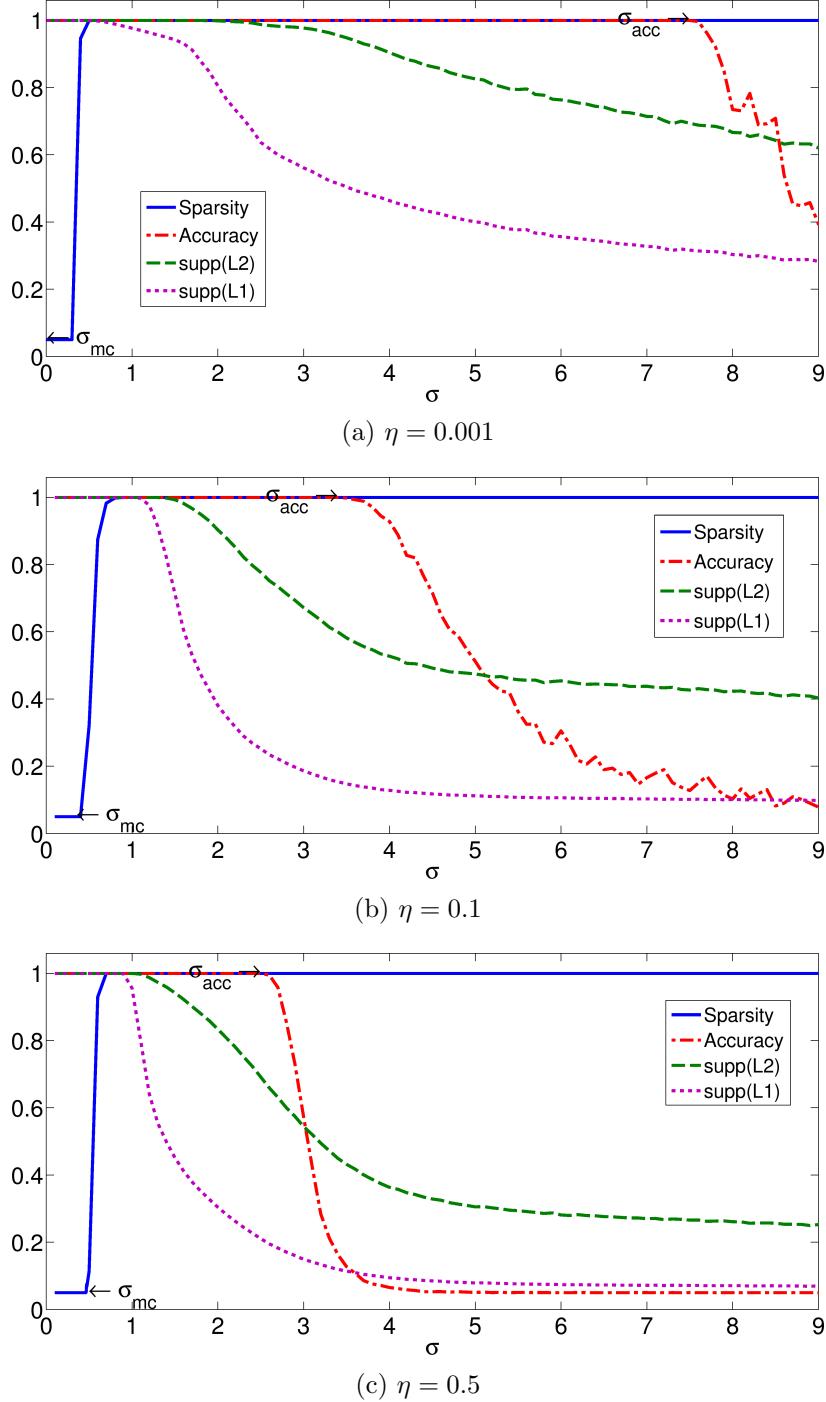


FIGURE 12.2. Average sparsity, accuracy,  $\text{supp}(\ell^2)$  and  $\text{supp}(\ell^1)$  (over 100 trials) as  $\sigma$  increased in the kernel setup. The annotations “ $\sigma_{\text{mc}}$ ” and “ $\sigma_{\text{acc}}$ ” denote the maximum  $\sigma$  for which Eq. (9.4) holds and for which maximum accuracy is obtained in Kernel SRC, respectively.

## 12.5. KEY FINDINGS

---

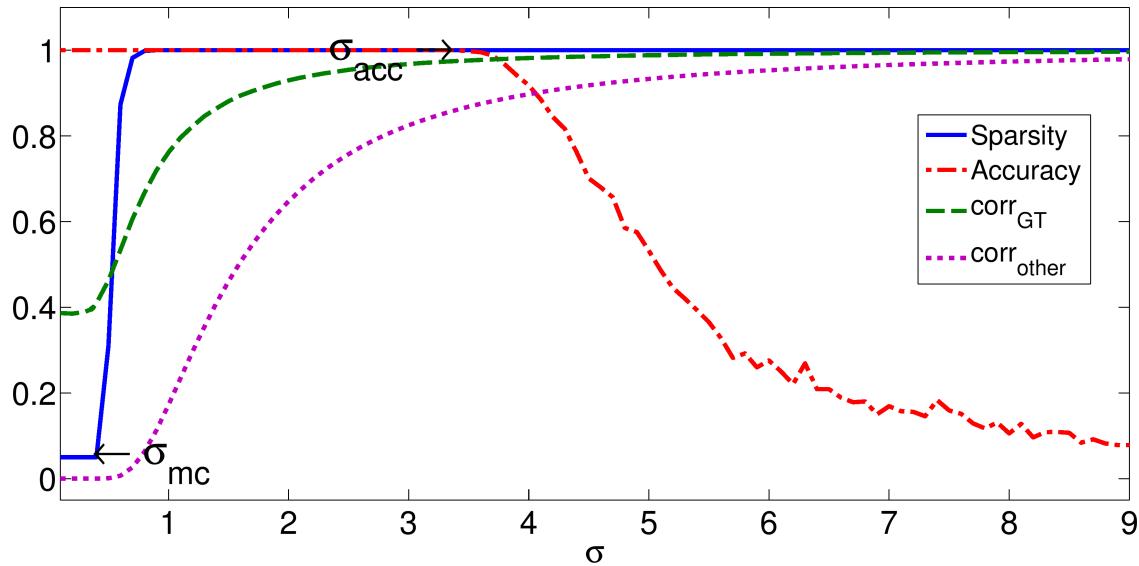


FIGURE 12.3. Median correlation (averaged over 100 trials) between the test sample  $\phi(\mathbf{y})$  and training samples in the same class ( $\text{corr}_{\text{GT}}$ ) and training samples in different classes ( $\text{corr}_{\text{other}}$ ) for the synthetic database with  $\eta = 0.1$ . Sparsity and accuracy are also displayed for comparison. Notice that the drop in accuracy occurs well after the jump in the correlation terms and sparsity.

## 12.5. KEY FINDINGS

---

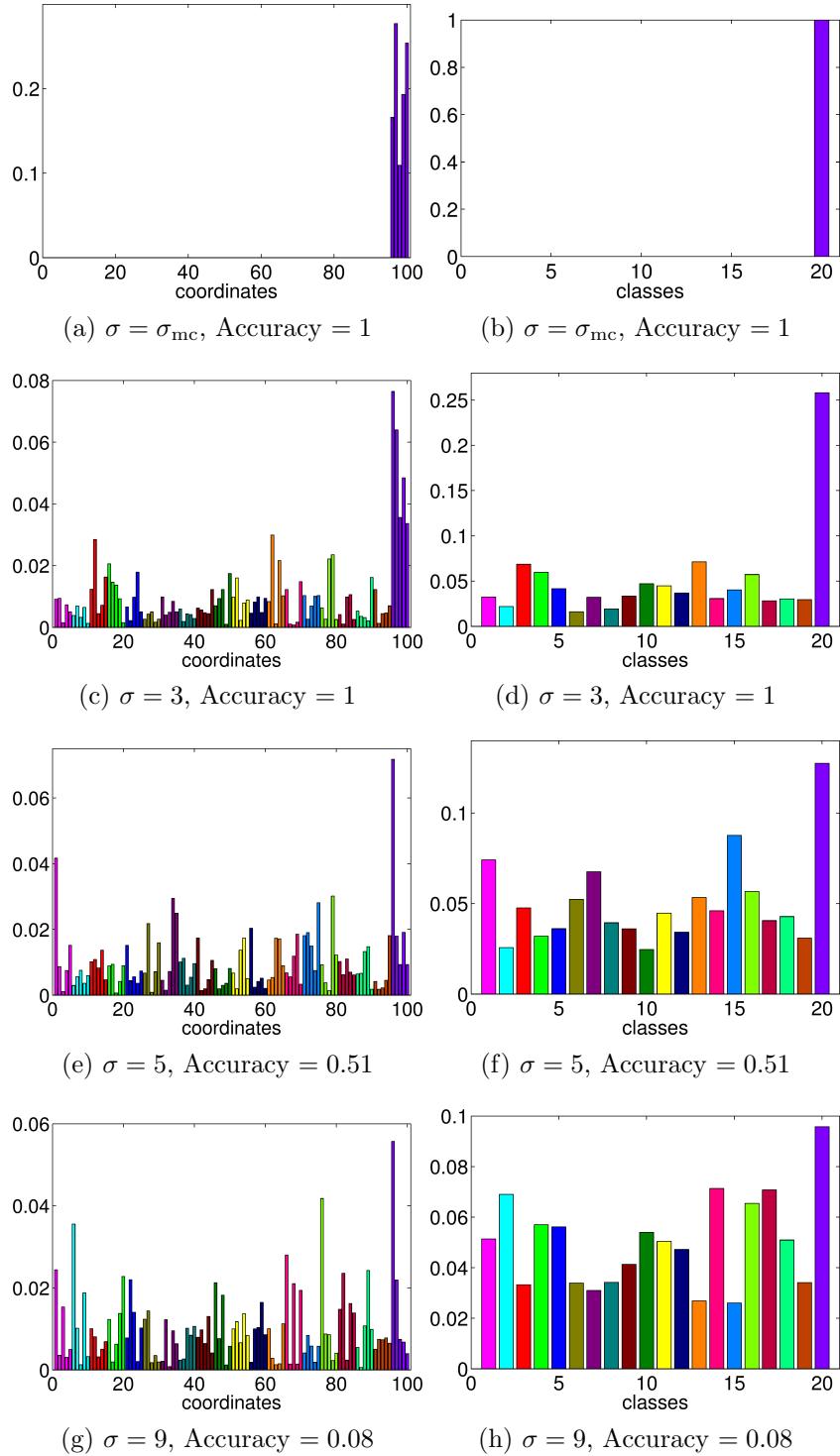


FIGURE 12.4. Average class contributions (over 100 trials) of coefficient vectors corresponding to class  $l = 20$  test samples. The colors denote the classes.

## CHAPTER 13

### Conclusion of Part 2

In Part 2 of this dissertation, we investigated the applicability of  $\ell^1/\ell^0$ -equivalence guarantees on dictionaries containing training samples. We detailed the inherent conflict between tightly-clustered classes—desirable for good classification—and the sufficient incoherence required by equivalence guarantees such as those based on restricted isometry and mutual coherence. In particular, we proved that under the assumptions of SRC, i.e., that class manifolds are linear subspaces spanned by their respective training data, Donoho et al.’s mutual coherence guarantees can only hold in the case that we have *exactly* enough training samples to span each lower-dimensional subspace. Considering that the performance of SRC generally improves as the training class size increases (recall that a primary goal of the LPCA-SRC algorithm from Part 1 was to enlarge the size of the training set in SRC), it is likely counter-productive for classification purposes to restrict the training set in this way. Further, despite the methods referenced in Part 1 to estimate the class manifold dimension, it is impractical to assume that such approaches will always work perfectly.

Despite not being able to prove  $\ell^1/\ell^0$ -equivalence on most class-structured data, we saw that it can indeed be achieved in some specific cases. Inspired by the random model of Wright and Ma to generate face image-like databases, we designed an experiment to test the ability of  $\ell^1$ -minimization to recover the sparsest solution on highly-correlated data. The results were mostly positive. We observed that in all cases,  $\ell^1$ -minimization recovered a solution closely approximating the sparsest solution (defined by generating the test sample as a linear combination of training samples in its ground truth class). Further, within-class correlation actually improved recovery relative to uniformly-random data, provided that the between-class correlation was sufficiently low, i.e., that the classes were sufficiently separated. In many cases,  $\ell^1$ -minimization exactly recovered the sparsest solution. Additionally, in the case that noise was added to the test sample, the correct support was found in nearly every case in which correlation was introduced.

---

We also considered the role of sparsity in the context of SRC and similar classification algorithms. One obstacle in determining this relationship is obtaining access to the sparsest solution for comparison without the aid of  $\ell^1/\ell^0$ -equivalence guarantees. Towards resolving this problem, we designed a nonlinear transform, based on kernel methods using the Gaussian kernel, to decrease the within-class mutual coherence while still maintaining class structure so that (hypothetically) provable equivalence and good classification could be simultaneously achieved. However, we found that the degree to which we had to decrease coherence in this setup meant that the test sample was no longer in the span of the training data, and so we were forced to limit our analysis to test samples artificially generated as linear combinations of their ground truth class training samples, as in Chapter 11. Though this to some extent limited the applicability of our experiment, the results clearly indicate that strict sparsity is not necessary for good classification in SRC. Instead, its success lies in its ability to correctly differentiate the coefficient magnitudes of training samples in different classes, i.e., to find *approximately* or *relatively* sparse solutions, in the case that the linear subspace assumption is observed and the classes themselves are not too correlated, i.e., not close together.

There is certainly much work to be done to quantify these findings. We mention two potential next steps: Eldar and Kupfinger's notion of *block-coherence* [36], with blocks corresponding to classes of the training database, might serve to make precise the meaning of between-class correlation; note that this was observed to play a role in both  $\ell^1/\ell^0$ -equivalence on highly-correlated data and SRC's classification performance. Additionally, the accuracy threshold detected in Chapter 12 might be better understood in the context of Wang et al.'s interpretation of SRC as a maximum margin-based classifier [98]. As an alternative to the thresholding route as suggested in Chapter 12, their work could be very helpful in rigorously defining the concept of *approximate sparsity* as it relates to the classification performance of SRC.

## CHAPTER 14

### Final Remarks

In this last chapter, we make some final comments regarding the relationship between classification and sparsity.

We saw on the sinusoidal waveform synthetic database in Chapter 5 that the sparsity enforced via  $\ell^1$ -minimization resulted in an extremely local approximation of the test sample, leading to good classification performance that improved as the sampling density of the class manifolds increased. In this case, *sparsity implied locality*, and as we witnessed via the  $\ell^2$ -regularized LDSL-SRC results, the converse was not true, at least to the extent necessary for good classification performance on this database.

However, the  $\ell^1$ -minimized solution does not always result in a local approximation of the test sample: On the ORL database, for instance, we saw that the coefficient vector contained significant coefficients from nonlocal training samples, as evidenced by the increased accuracy of SRC<sub>pruned</sub> over SRC. This demonstrates that enforcing locality via a simple dictionary pruning step can improve classification performance over  $\ell^1$ -minimization alone.

These examples illustrate that neither sparsity nor locality alone is enough to produce good classification results in all situations. The algorithm LPCA-SRC is a step in combining both properties into a single classifier; however, there is certainly room for improvement in terms of computational complexity and robustness to data noise because of the use of tangent vectors. There is also more work to be done in terms of the easy estimation of tangent vector error and providing the user with concrete recommendations and ideally, performance guarantees.

Given the findings in Part 2, we are confident that  $\ell^1$ -minimization, even on highly-correlated data, will find a sparse or relatively-sparse coefficient vector provided that the test sample can be approximated reasonably well using samples from a single class and that the training set contains enough class structure to make the classification problem viable. This is despite the fact that the mutual coherence  $\ell^1/\ell^0$ -equivalence guarantees—the only tractable approach discussed in Chapter

---

9—provably cannot be applied in general, provided that the assumptions in SRC are met. We observed that, though the exact sparsest solution may be found using  $\ell^1$ -minimization provided that the classes are well-separated, this may be unnecessary, due to the robust classification mechanism in SRC. Further, the relatively sparse solutions awarded by  $\ell^1$ -minimization can lead to significantly better classification results than the dense solutions obtained using  $\ell^2$ -regularization, particularly when the feature dimension is small or, as our experiments in Chapter 7 suggest, in the case that multiple good per-class approximations of the test sample exist, as on the sinusoidal waveform synthetic database. More work is needed to gain deeper understanding of these relationships.

Our results are very much dependent on the representation-based approach to classification being used, and it is important to note that this method is not applicable on some data sets. For example, modeling a test sample as an approximate linear combination of its same-class training samples is not appropriate when the training data do not lie on smooth manifolds (or not on manifolds at all) or when there is only a single training sample per class.<sup>1</sup> Further, we have focused very little on the fact that these methods require the training samples to be normalized; when essential class-structure is contained in the norms of the samples, these methods fail.<sup>2</sup>

We end this dissertation by mentioning the *bet on sparsity* principle of Hastie et al. [48]. In the context of variable selection for regression and classification problems, this principle states that we should use a procedure that will find sparse solutions (i.e., we should assume a sparseness prior), because *no procedure will do well if the solution is dense* [48]. That is, consider the dictionary with columns representing variables (or features) and rows representing samples. (Recall that the opposite is true in the representation-based classification scenario.) Given that there are more variables than samples so that the resulting system is underdetermined, the goal is to express a vector of responses  $\mathbf{y}$  as an approximate linear combination of these variables. If the (ground truth) coefficients of this approximation are not sparse, then we do not have enough samples to correctly identify them [48]. For example, both  $\ell^1$ -regularization and the denser  $\ell^2$ -regularization approaches will fail to recover this solution. So we might as well assume that the coefficients are sparse and move forward.

---

<sup>1</sup>However, SRC has been extended to the single sample face recognition scenario. See Deng et al.’s algorithm [26].

<sup>2</sup>Interestingly, the kernel trick has been used to address this problem. See Zhang et al.’s paper [116].

---

Let us detail an analogous principle in the representation-based classification scenario. Given that we wish to express a test sample  $\mathbf{y}$  over a training set  $X_{\text{tr}} \in \mathbb{R}^{m \times N_{\text{tr}}}$  with  $m \ll N_{\text{tr}}$ , then if  $\mathbf{y}$  is not an approximate linear combination of its same-class training samples, it is unclear how we will be able to determine the correct class using any representation-based classifier. On the other hand, if a representation using the correct class exists and is not sparse, i.e., if  $L$  is small, then either the sparsest representation will consist of a subset of these class coefficients (in which case this class will be identified), or the sparsest representation will contain nonzero coefficients corresponding to training samples from several cases. If this happens, then there is no reason to think (in general) that an incorrect class will contribute more to this representation than the correct one (barring outliers and heavy sample corruption, of course). Even in the case that  $\mathbf{y}$  lies in the intersection of two class manifolds, we will at worst be aware of the uncertainty of the classification decision via the similarity in class residuals  $\text{err}_l(\mathbf{y})$  and have identified the two most likely classes. Note that, in this scenario, the block-sparse SSR methods would do no better, and CRC-RLS may do substantially worse.

Thus, in considering representation-based classification, one should “bet on sparsity” via  $\ell^1$ -minimization. Even in the case that the class approximation is not found, it is unlikely that any other approach to finding it will do better.

## APPENDIX A

# Proof of Tangent Bound Modification

The main result in Kaslovsky and Meyer’s paper [57] bounds the error between an approximated tangent plane computed using standard local PCA and the true tangent plane. We refer to this as the “tangent bound.” One can view the offline portion of the classification algorithm LPCA-SRC as a method of generating new training samples in the form of shifted and scaled basis vectors of the approximated tangent plane; we call these basis elements “tangent vectors.” Notably, LPCA-SRC does not use standard local PCA to compute the tangent vectors but instead the local PCA technique of Singer and Wu [87]. Thus, in order to use the tangent bound to estimate the error in the newly-generated training samples (i.e., the tangent vectors), as was done in Chapter 6 above, the bound in the Kaslovsky and Meyer paper [57] must be amended to account for the difference in local PCA implementations. We do this here.

### A.1. Tangent Bound Details

**A.1.1. Geometric Setup.** Given a smooth,  $d$ -dimensional manifold  $\mathcal{M} \subseteq \mathbb{R}^m$  and points randomly sampled from  $\mathcal{M}$ , we can approximate the tangent (hyper)plane to the manifold at a sampled point  $\mathbf{x}_0$ . Local PCA is a standard tool to do this. To perform standard local PCA, we determine the  $n$  nearest neighbors of  $\mathbf{x}_0$  from the sampled data set and compute the covariance matrix from these samples. The first  $d$  eigenvectors of the covariance matrix are then used as the basis vectors of the approximated tangent plane.

Kaslovsky and Meyer’s tangent bound [57] determines (with high probability) the maximum angle between the true tangent plane and the approximated tangent plane. An important part of the authors’ theory is a specific *geometric setup*, involving a (theoretical) shift and rotation of the manifold under which this angle is invariant. The geometric setup identifies  $\mathbf{x}_0$  with the origin and aligns the principal axes of curvature of  $\mathcal{M}$  at  $\mathbf{x}_0$  with the coordinate axes. This allows for the separation of the linear and curvature components of each of the  $n$  neighbors of  $\mathbf{x}_0$ . That is, a

## A.1. TANGENT BOUND DETAILS

---

nearby point  $\mathbf{x}$  in the shifted and rotated coordinate system can be written as

$$\mathbf{x} = \boldsymbol{\ell} + \mathbf{c} + \mathbf{e} = \begin{bmatrix} \ell_1 \\ \vdots \\ \ell_d \\ 0 \\ \vdots \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ \vdots \\ 0 \\ c_{d+1} \\ \vdots \\ c_m \end{bmatrix} + \begin{bmatrix} e_1 \\ \vdots \\ e_d \\ e_{d+1} \\ \vdots \\ e_m \end{bmatrix} \in \mathbb{R}^m,$$

where the vectors  $\boldsymbol{\ell}$ ,  $\mathbf{c}$ , and  $\mathbf{e}$  refer to the linear, curvature, and noise components of  $\mathbf{x}$ , respectively.

If the  $n$  nearest neighbors of  $\mathbf{x}_0$  in the data set obey certain sampling assumptions, i.e., if their linear components are assumed to be distributed uniformly within a neighborhood of  $\mathbf{x}_0$  on the tangent plane at  $\mathbf{x}_0$ , then we can model this situation as follows: For  $\mathbf{x} = \boldsymbol{\ell} + \mathbf{c} + \mathbf{e}$  in the geometric setup, the first  $d$  coordinates of  $\boldsymbol{\ell} \in \mathbb{R}^m$  can be viewed as random samples from the  $d$ -dimensional tangent plane at  $\mathbf{x}_0$ , i.e., the span of the first  $d$  axes, uniform within  $B_0^d(r)$ , the  $d$ -dimensional ball of radius  $r$  centered at  $\mathbf{x}_0 = \mathbf{0}$ . Here,  $r$  is used to define the local neighborhood, and we stress that  $r$  bounds the distance between  $\mathbf{x}_0$  and  $\mathbf{x}$  as measured *along the tangent plane*, i.e.,  $\|\boldsymbol{\ell}\| \leq r$ .

The curvature component  $\mathbf{c} \in \mathbb{R}^m$  can be viewed as having nonzero coordinates defined by

$$(A.1) \quad c_i := \frac{1}{2} \left( \kappa_1^{(i)} \ell_1^2 + \cdots + \kappa_d^{(i)} \ell_d^2 \right), \quad i = d+1, \dots, m,$$

where  $\{\kappa_1^{(i)}, \dots, \kappa_d^{(i)}\}$  are the principal curvatures of  $\mathcal{M}$  at  $\mathbf{x}_0$  in the  $i$ th normal direction,  $i = d+1, \dots, m$ .<sup>1</sup> Lastly, the noise component  $\mathbf{e} \in \mathbb{R}^m$  can be viewed as being drawn from  $\mathcal{N}(\mathbf{0}, \sigma^2 I_m)$  (assuming that this is an appropriate characterization of the sample noise).

**A.1.2. Notation.** We assume throughout this discussion that local data samples can be decomposed as above, i.e., we view the samples as points in the shifted and rotated coordinate system, generated as above. Given the local samples  $\mathbf{x}^{(k)} = \boldsymbol{\ell}^{(k)} + \mathbf{c}^{(k)} + \mathbf{e}^{(k)} \in \mathbb{R}^m$ ,  $1 \leq k \leq n$ , define the

---

<sup>1</sup>The reader may recall from the tangent distance theorem in Chapter 6 that Eq. (A.1) is actually the second-order Taylor expansion of the  $i$ th coordinate of the curvature component. Citing the numerical results of Tyagi et al. [94], Kaslovsky and Meyer claim that this expansion is sufficiently general in this context [57].

## A.1. TANGENT BOUND DETAILS

---

matrix  $X := [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}] \in \mathbb{R}^{m \times n}$ . Then we can write

$$X = L + C + E,$$

where  $L := [\ell^{(1)}, \dots, \ell^{(n)}]$  is the linear/tangent component,  $C := [c^{(1)}, \dots, c^{(n)}]$  is the curvature component, and  $E := [e^{(1)}, \dots, e^{(n)}]$  is the noise component. To approximate the tangent plane, we aim to recover  $L$  from  $X$ . As mentioned, standard local PCA's solution to this problem is to compute the first  $d$  eigenvectors of the covariance matrix of  $X$  and set these to be the approximated tangent plane basis vectors. Kaslovsky and Meyer's tangent bound compares this solution to the true tangent plane.

Let  $\tilde{X} := X - (1/n) \sum_{k=1}^n \mathbf{x}^{(k)}$  be the *centered* version of  $X$ . Kaslovsky and Meyer define the following eigendecompositions [57]:

$$\frac{1}{n} \tilde{X} \tilde{X}^\top = \hat{U} \hat{\Lambda} \hat{U}^\top = \begin{bmatrix} \hat{U}_1 & \hat{U}_2 \end{bmatrix} \begin{bmatrix} \hat{\Lambda}_1 & 0 \\ 0 & \hat{\Lambda}_2 \end{bmatrix} \begin{bmatrix} \hat{U}_1 & \hat{U}_2 \end{bmatrix}^\top$$

and

$$\frac{1}{n} \tilde{L} \tilde{L}^\top = U \Lambda U^\top = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \Lambda_1 & 0 \\ 0 & \Lambda_2 \end{bmatrix} \begin{bmatrix} U_1 & U_2 \end{bmatrix}^\top,$$

where  $\hat{U}_1$  and  $U_1$  contain the first  $d$  columns of  $\hat{U}$  and  $U$ , respectively, and the diagonal matrices of eigenvalues,  $\hat{\Lambda}$  and  $\Lambda$ , are similarly divided. It is assumed that the eigenvalues are in decreasing order. It follows that the angle between the true and approximated tangent planes can be measured by computing the Frobenius norm of the difference between their respective orthogonal projectors:

$$(A.2) \quad \|\hat{U}_1 \hat{U}_1^\top - U_1 U_1^\top\|_F.$$

Note that  $\hat{U}_1$  contains the approximated tangent plane basis vectors found by standard local PCA.

To bound  $\|\hat{U}_1 \hat{U}_1^\top - U_1 U_1^\top\|_F$ , Kaslovsky and Meyer [57] expand the covariance matrix of  $X$  as

$$\frac{1}{n} \tilde{X} \tilde{X}^\top = (\tilde{L} + \tilde{C} + \tilde{E})(\tilde{L} + \tilde{C} + \tilde{E})^\top = \tilde{L} \tilde{L}^\top + \Delta,$$

## A.2. SINGER AND WU'S LOCAL PCA

---

where

$$(A.3) \quad \Delta := \frac{1}{n}(\tilde{C}\tilde{C}^\top + \tilde{E}\tilde{E}^\top + \tilde{L}\tilde{C}^\top + \tilde{C}\tilde{L}^\top + \tilde{L}\tilde{E}^\top + \tilde{E}\tilde{L}^\top + \tilde{C}\tilde{E}^\top + \tilde{E}\tilde{C}^\top),$$

and  $\tilde{L}$ ,  $\tilde{C}$ , and  $\tilde{E}$  denote the centered versions of  $L$ ,  $C$ , and  $E$ , respectively. The authors next compute a lower bound for  $\lambda_d$ , the  $d$ th largest eigenvalue of  $(1/n)\tilde{L}\tilde{L}^\top$  (equivalently, the  $(d,d)$ th entry of  $\Lambda_1$ ), and upper bounds for the matrix norms  $\|U_1^\top \Delta U_1\|_F$ ,  $\|U_2^\top \Delta U_2\|_F$ ,  $\|U_2^\top \Delta U_1\|_F$ , and  $\|U_1^\top \Delta U_2\|_F$ . They then use these quantities in the following theorem:

**THEOREM A.1.1** (Davis and Kahan [25], Stewart [89]). *Set*

$$\delta := \lambda_d - \|U_1^\top \Delta U_1\|_F - \|U_2^\top \Delta U_2\|_F.$$

*Define the following conditions:*

- *Condition 1:*  $\delta > 0$ ,
- *Condition 2:*  $\|U_1^\top \Delta U_2\|_F \|U_2^\top \Delta U_1\|_F < \frac{1}{4}\delta^2$ .

*If both conditions hold, then*

$$\|\hat{U}_1 \hat{U}_1^\top - U_1 U_1^\top\|_F \leq 2\sqrt{2} \frac{\|U_2^\top \Delta U_1\|_F}{\delta}.$$

**A.1.3. Sampling Assumptions.** Kaslovsky and Meyer [57] assume that the number of neighbors  $n$  satisfies  $n \geq O(d \log d)$ . More generally, the spirit of their analysis is based on a generous sampling density, something to keep in mind when considering the applicability of (our modification to) their theorem in the context of classification problems.

## A.2. Singer and Wu's Local PCA

In standard local PCA, the approximated tangent plane basis vectors are found by computing the eigenvectors of the covariance matrix  $(1/n)\tilde{X}\tilde{X}^\top$ . However, the local PCA technique of Singer and Wu [87] computes the eigenvectors of a slightly different covariance matrix. Given the local sample matrix  $X = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]$ , we define

$$Y := (X - [\mathbf{x}_0, \dots, \mathbf{x}_0])D_0,$$

### A.3. MODIFYING KASLOVSKY AND MEYER'S SETUP

---

for  $\mathbf{x}_0$  the point on the manifold at which we are computing the tangent plane.  $D_0$  is a diagonal matrix whose  $(j,j)$ th entry contains an increasing function of the inverse distance between  $\mathbf{x}_0$  and  $\mathbf{x}^{(j)}$ ,  $1 \leq j \leq n$ , so that points nearer to  $\mathbf{x}_0$  contribute more to the construction of the tangent plane at  $\mathbf{x}_0$ . The approximated tangent plane basis vectors in Singer and Wu's local PCA are given by the eigenvectors of the matrix  $YY^\top$ , equivalently, the left singular vectors of  $Y$ . Consider the difference between this method and standard local PCA, which computes the eigenvectors of  $(1/n)\tilde{X}\tilde{X}^\top$ , where  $\tilde{X} = (X - [\boldsymbol{\mu}, \dots, \boldsymbol{\mu}])I$  and  $\boldsymbol{\mu}$  is the mean column of  $X$  given by  $\boldsymbol{\mu} := (1/n) \sum_{k=1}^n \mathbf{x}^{(k)}$ .

### A.3. Modifying Kaslovsky and Meyer's Setup

**A.3.1. Amending the Geometric Setup.** What effect does the shift by  $\mathbf{x}_0$  in Singer and Wu's local PCA, in contrast to the shift by  $\boldsymbol{\mu}$  in standard local PCA, have on the geometric setup and the model it produces? Since the geometric setup automatically (i.e., *inherently*) shifts the points by  $\mathbf{x}_0$  in order to treat this point as the origin, it suffices to consider the *uncentered* matrix  $X$  in the tangent bound analysis. Thus the Singer and Wu case with  $D_0 = I$  is simply Kaslovsky and Meyer's original setup with  $X$  replacing the centered matrix  $\tilde{X}$ . In other words, we only need to consider  $Y = XD_0$  in the geometric setup and do not need to make any explicit compensation for the centering by  $\mathbf{x}_0$ .

Next we consider the effect of the weight matrix  $D_0$ . Let  $D_0$  have diagonal entries  $d_1, \dots, d_n$ . We expand

$$Y = (X - [\mathbf{x}_0, \dots, \mathbf{x}_0])D_0 = [d_1(\mathbf{x}^{(1)} - \mathbf{x}_0), \dots, d_n(\mathbf{x}^{(n)} - \mathbf{x}_0)].$$

Thus when  $\mathbf{x}_0$  is treated as the local origin, the points we must consider are the scaled samples

$$\{d_1\mathbf{x}^{(1)}, \dots, d_n\mathbf{x}^{(n)}\}.$$

Recall that in local PCA in the LPCA-SRC algorithm (i.e., Algorithm 4), the samples in  $X$  are sorted in increasing order of distance to  $\mathbf{x}_0$ . However, this is purely for convenience: we will obtain the same eigenvectors regardless of the arrangement of the columns of  $X$ . Thus  $d_k$  should refer to the weight of the  $k$ th sample  $\mathbf{x}^{(k)}$ , which may not necessarily be the  $k$ th closest sample to  $\mathbf{x}_0$ .

### A.3. MODIFYING KASLOVSKY AND MEYER'S SETUP

---

Next we consider how to model (or view as randomly-generated samples) the points  $d_1 \mathbf{x}^{(1)}, \dots, d_n \mathbf{x}^{(n)}$  in the geometric setup. We have

$$d_k \mathbf{x}^{(k)} = d_k(\boldsymbol{\ell}^{(k)} + \mathbf{c}^{(k)} + \mathbf{e}^{(k)}) = d_k \boldsymbol{\ell}^{(k)} + d_k \mathbf{c}^{(k)} + d_k \mathbf{e}^{(k)},$$

for  $1 \leq k \leq n$ . That is, the linear, curvature, and noise components are simply scaled by  $d_k$ . However, the (weighted) linear components  $d_k \boldsymbol{\ell}^{(k)}$  are no longer uniformly distributed within  $B_{\mathbf{0}}^d(\mathbf{x}_0)$ . Of course, this distribution depends on  $d_k$ .

In the LPCA-SRC algorithm, we set

$$\begin{aligned} d_k &:= \sqrt{1 - \frac{\|\mathbf{x}^{(k)}\|^2}{s}} \\ &= \sqrt{1 - \frac{\|\boldsymbol{\ell}^{(k)}\|^2 + \|\mathbf{c}^{(k)}\|^2 + \|\mathbf{e}^{(k)}\|^2 + 2 \langle \boldsymbol{\ell}^{(k)} + \mathbf{c}^{(k)}, \mathbf{e}^{(k)} \rangle}{s}}, \end{aligned}$$

where  $s$  is a scaling factor satisfying  $s > \max_{1 \leq k \leq n} \|\mathbf{x}^{(k)}\|^2$ . The second equality follows from

$$\langle \boldsymbol{\ell}^{(k)}, \mathbf{e}^{(k)} \rangle = 0.$$

**A.3.2. Amending the Notation.** We write  $XD_0 = LD_0 + CD_0 + ED_0$  and

$$\frac{1}{n} XD_0(XD_0)^T = \frac{1}{n} (LD_0 + CD_0 + ED_0)(LD_0 + CD_0 + ED_0)^T = \frac{1}{n} LD_0(LD_0)^T + \Theta,$$

where

$$\begin{aligned} \Theta &:= \frac{1}{n} (CD_0(CD_0)^T + ED_0(ED_0)^T + LD_0(CD_0)^T + CD_0(LD_0)^T \\ (A.4) \quad &\quad + LD_0(ED_0)^T + ED_0(LD_0)^T + CD_0(ED_0)^T + CD_0(CD_0)^T). \end{aligned}$$

Similarly to Kaslovsky and Meyer's notation, we define the terms in the eigendecomposition of  $(1/n)XD_0(XD_0)^T$  using

$$(A.5) \quad \frac{1}{n} XD_0(XD_0)^T = \begin{bmatrix} \widehat{W}_1 & \widehat{W}_2 \end{bmatrix} \begin{bmatrix} \widehat{\Lambda}_1^W & 0 \\ 0 & \widehat{\Lambda}_2^W \end{bmatrix} \begin{bmatrix} \widehat{W}_1 & \widehat{W}_2 \end{bmatrix}^T.$$

#### A.4. SCALING FACTOR

---

and

$$(A.6) \quad \frac{1}{n} LD_0(LD_0)^T = \begin{bmatrix} W_1 & W_2 \end{bmatrix} \begin{bmatrix} \Lambda_1^W & 0 \\ 0 & \Lambda_2^W \end{bmatrix} \begin{bmatrix} W_1 & W_2 \end{bmatrix}^T.$$

Note that we have replaced the  $U$ 's with  $W$ 's,  $\Lambda$ 's with  $\Lambda^W$ 's, and  $\Delta$  with  $\Theta$  to differentiate the Singer and Wu case from the authors' original analysis [57].

Our goal is to bound

$$(A.7) \quad \|\widehat{W}_1 \widehat{W}_1^T - W_1 W_1^T\|_F$$

by modifying Kaslovsky and Meyer's bounds and inserting them into Theorem A.1.1. Note that Theorem A.1.1 does not require any special structure of the input matrix (from which the eigen-decomposition stems), e.g., this matrix does not need to have columns centered around  $\mathbf{0}$ .

As a last note, recall that Singer and Wu's local PCA uses the eigenvectors of  $XD_0(XD_0)^T$ , not those of  $(1/n)XD_0(XD_0)^T$ . However, since  $XD_0(XD_0)^T = LD_0(LD_0)^T + n\Theta$ , it suffices to simply multiply the relevant terms in Theorem A.1.1 by  $n$  in order to account for this difference. We do so after all the bounds are obtained.

#### A.4. Scaling Factor

In LPCA-SRC, we set the scaling factor  $s := \|\boldsymbol{x}^{(n+1)} - \boldsymbol{x}_0\|^2$ , where  $\boldsymbol{x}^{(n+1)}$  is defined to be the farthest away of the  $n + 1$  samples (the  $(n + 1)$ -nearest neighbors of  $\boldsymbol{x}_0$  from the same class) from  $\boldsymbol{x}_0$ . Equivalently,  $\boldsymbol{x}^{(n+1)}$  is the sample with the largest norm in the geometric setup, and thus it depends on all of the samples  $\boldsymbol{x}^{(k)}$ ,  $1 \leq k \leq n$ . This choice of  $s$  is somewhat arbitrary; it provides a simple way of defining the neighborhood radius used in Singer and Wu's local PCA [87]. But all that is necessary in their algorithm is that  $s$  satisfy  $s > \|\boldsymbol{x}^{(k)} - \boldsymbol{x}_0\|^2$  for  $1 \leq k \leq n$ . Further, without being able to assume that  $s$  is independent of  $X$ , it will be very difficult to amend Kaslovsky and Meyer's analysis to the Singer and Wu case. Thus we define  $s$  slightly differently here.

Treating  $\boldsymbol{x}_0$  as the origin, we have that

$$\|\boldsymbol{x}^{(k)} - \mathbf{0}\|^2 = \|\boldsymbol{\ell}^{(k)}\|^2 + \|\boldsymbol{c}^{(k)}\|^2 + \|\boldsymbol{e}^{(k)}\|^2 + 2 \langle \boldsymbol{\ell}^{(k)} + \boldsymbol{c}^{(k)}, \boldsymbol{e}^{(k)} \rangle$$

#### A.4. SCALING FACTOR

---

$$\begin{aligned} &\leq r^2 + \frac{(K^{(+)})^2 r^4}{4} + \sum_{j=1}^m (e_i^{(k)})^2 + 2r \sum_{j=1}^d |e_j^{(k)}| + K^{(+)} r^2 \left( \sum_{i=d+1}^m (e_i^{(k)})^2 \right)^{1/2} \\ &\leq \frac{(K^{(+)})^2 r^4}{4} + K^{(+)} r^2 \sqrt{m-d} (C\sigma) + r^2 + m(C\sigma)^2 + 2rd(C\sigma) \end{aligned}$$

for all samples  $\mathbf{x}^{(k)}$ ,  $1 \leq k \leq n$ . Here,  $K^{(+)}$  and  $K_i^{(+)}$  are Kaslovsky and Meyer's curvature constants defined using

$$(A.8) \quad K_i^{(+)} := \left( \sum_{j=1}^d |\kappa_j^{(i)}|^2 \right)^{1/2} \text{ and } K^{(+)} := \left( \sum_{i=d+1}^m (K_i^{(+)})^2 \right)^{1/2}$$

(we derive the bound  $\|\mathbf{c}\|^2 \leq (K^{(+)})^2 r^4 / 4$  in Section A.5.4 below). Lastly,  $C \in \mathbb{N}$  is chosen so that  $|e_i^{(k)}| \leq \sigma$  with the desired probability. For example, setting  $C = 3$  implies that the above bound holds with probability at least 99.7%.

Thus, in order to satisfy  $s > \|\mathbf{x}^{(k)} - \mathbf{x}_0\|^2$ , we set

$$(A.9) \quad s := \frac{(K^{(+)})^2 r^4}{4} + K^{(+)} r^2 \sqrt{m-d} (C\sigma) + r^2 + m(C\sigma)^2 + 2rd(C\sigma) + \epsilon,$$

for  $\epsilon$  a small, positive constant.

**REMARK A.4.1.** We briefly note that this bound on  $\|\mathbf{x}^{(k)}\|^2$ ,  $1 \leq k \leq n$ , could be large, and so this assignment for the scaling factor  $s$  could offer little differentiation in the weights  $d_k$  (making them all close to 1) and rendering this important aspect of Singer and Wu's local PCA ineffective. At the end of our analysis, we consider what happens when a smaller value of  $s$  is used. See Section A.7.

**A.4.1. Breakdown of Terms.** Following Kaslovsky and Meyer's lead, we consider the terms that need to be bounded. By Theorem A.1.1 (accounting for the change in notation), it is sufficient to bound the terms  $\|W_1^\top \Theta W_1\|_F$ ,  $\|W_2^\top \Theta W_2\|_F$ ,  $\|W_2^\top \Theta W_1\|_F$ ,  $\|W_1^\top \Theta W_2\|_F$ , and  $\lambda_d^W$ , the  $d$ th largest eigenvalue of  $(1/n)L D_0 (L D_0)^\top$ , in order to bound Eq. (A.7). In particular, the bound on  $\lambda_d^W$  must be a lower bound.

Because  $\Theta$  is symmetric (like Kaslovsky and Meyer's  $\Delta$ ), we have that

$$\|W_1^\top \Theta W_2\|_F = \|W_2^\top \Theta W_1\|_F,$$

#### A.4. SCALING FACTOR

---

and so we consider only the subscript pairs  $(i, j) \in \{(1, 1), (2, 2), (2, 1)\}$ . Further, the equalities  $W_1^\top C = 0$  and  $W_2^\top L = 0$  follow from orthogonality, similarly to the standard local PCA case. Using these properties, the expansion of  $\Theta$  using Eq. (A.4) and the triangle inequality, the terms that need bounding are

- (1)  $\lambda_d^W$ ,
- (2)  $\left\| W_1^\top \left( \frac{1}{n} LD_0 (ED_0)^\top \right) W_1 \right\|_F$ ,
- (3)  $\left\| W_1^\top \left( \frac{1}{n} ED_0 (ED_0)^\top \right) W_1 \right\|_F$ ,
- (4)  $\left\| W_2^\top \left( \frac{1}{n} CD_0 (CD_0)^\top \right) W_2 \right\|_F$ ,
- (5)  $\left\| W_2^\top \left( \frac{1}{n} CD_0 (ED_0)^\top \right) W_2 \right\|_F$ ,
- (6)  $\left\| W_2^\top \left( \frac{1}{n} ED_0 (ED_0)^\top \right) W_2 \right\|_F$ ,
- (7)  $\left\| W_2^\top \left( \frac{1}{n} CD_0 (LD_0)^\top \right) W_1 \right\|_F$ ,
- (8)  $\left\| W_2^\top \left( \frac{1}{n} ED_0 (LD_0)^\top \right) W_1 \right\|_F$ ,
- (9)  $\left\| W_2^\top \left( \frac{1}{n} CD_0 (ED_0)^\top \right) W_1 \right\|_F$ ,
- (10)  $\left\| W_2^\top \left( \frac{1}{n} ED_0 (ED_0)^\top \right) W_1 \right\|_F$ .

Additionally, some of these Frobenius norm bounds will utilize upper bounds on the eigenvalues of the matrices  $(1/n)LD_0(LD_0)^\top$ ,  $(1/n)CD_0(CD_0)^\top$ , and  $(1/n)ED_0(ED_0)^\top$ . These are determined first.

**REMARK A.4.2.** *As implied in this detailed introduction, our derived tangent bound in the case of Singer and Wu’s local PCA will be nearly entirely based on the derivation of Kaslovsky and Meyer’s tangent bound [57]. For example, the terms listed above are exactly analogous to the terms they bound in their work. In nearly all instances, the methods used to obtain the bounds in the Singer and Wu case are exactly those used in Kaslovsky and Meyer’s proofs, and they should be given full credit for these methods. To avoid redundancy, we will refrain from continuously writing “as the authors do...” or “analogously to Kaslovsky and Meyer...” and trust that the reader is well-aware that this is the case wherever applicable.*

*With this in mind, we will not “re-prove” Kaslovsky and Meyer’s careful (and lengthy) analysis here. Instead, we primarily focus on the amendments needed to handle the Singer and Wu case. The reader is referred to Kaslovsky and Meyer’s paper [57] for details and an instructive interpretation of their tangent bound.*

### A.5. Bounding Terms

**A.5.1. Linear Eigenvalues.** We will determine both upper and lower bounds for the eigenvalues of  $(1/n)LD_0(LD_0)^\top$  using the following theorem:

**THEOREM A.5.1** (Tropp [93]). *Let  $\{X_k\}$  be a finite sequence of independent, random, self-adjoint matrices satisfying  $X_k \succeq 0$  and  $\lambda_{\max}(X_k) \leq \lambda_\infty$  almost surely for all  $k$ , where  $\lambda_{\max}(X_k)$  denotes the largest eigenvalue of  $X_k$ . Set*

$$\mu_{\min} := \lambda_{\min}\left(\sum_{k=1}^n \mathbb{E}[X_k]\right), \quad \mu_{\max} := \lambda_{\max}\left(\sum_{k=1}^n \mathbb{E}[X_k]\right).$$

Then

$$\mathbb{P}\left[\lambda_{\min}\left(\sum_{k=1}^n X_k\right) \leq (1-\delta)\mu_{\min}\right] \leq d\left[\frac{e^{-\delta}}{(1-\delta)^{1-\delta}}\right]^{\frac{\mu_{\min}}{\lambda_\infty}}, \quad \text{for } \delta \in [0, 1]$$

and

$$\mathbb{P}\left[\lambda_{\max}\left(\sum_{k=1}^n X_k\right) \geq (1+\delta)\mu_{\max}\right] \leq d\left[\frac{e^{-\delta}}{(1+\delta)^{1+\delta}}\right]^{\frac{\mu_{\max}}{\lambda_\infty}}, \quad \text{for } \delta \geq 0.$$

Noting that

$$\frac{1}{n}LD_0(LD_0)^\top = \frac{1}{n} \sum_{k=1}^n d_k \ell^{(k)} (d_k \ell^{(k)})^\top = \frac{1}{n} \sum_{k=1}^n d_k^2 \ell^{(k)} (\ell^{(k)})^\top,$$

we set  $X_k := (1/n) d_k^2 \ell^{(k)} (\ell^{(k)})^\top$ . Clearly,  $X_k$  is self-adjoint and semi-positive definite. We need to determine the terms  $\lambda_\infty$ ,  $\mu_{\min}$ , and  $\mu_{\max}$ . Consider

$$\lambda_{\max}\left(\frac{1}{n} d_k^2 \ell^{(k)} \ell^{(k)\top}\right) = \frac{1}{n} \lambda_{\max}(d_k^2 \ell^{(k)} \ell^{(k)\top}) = \frac{1}{n} \|d_k \ell^{(k)}\|_2^2 \leq \frac{1}{n} \|\ell^{(k)}\|_2^2 \leq \frac{r^2}{n}.$$

Since this holds for all  $1 \leq k \leq n$ , we set  $\lambda_\infty := r^2/n$ .

Next, towards determining  $\mu_{\min}$  and  $\mu_{\max}$ , we compute  $\mathbb{E}[X_k]$ . Since the distribution of the entries of  $X_k$  is the same for all  $k$ , we can drop the “ $k$ ” notation, letting  $X$ ,  $\ell$ , and  $x$  denote the random variables from which the realizations  $X^{(k)}$ ,  $\ell^{(k)}$ , and  $x^{(k)}$ , respectively, are drawn. Similarly,

## A.5. BOUNDING TERMS

---

we use  $\omega$  to denote the random weight variable (from which  $d_k$  is drawn), i.e.,

$$(A.10) \quad \omega := \sqrt{1 - \frac{\|\mathbf{x}\|^2}{s}}.$$

It follows that we need to compute  $\mathbb{E}[X] = (1/n)\mathbb{E}[\omega^2 \ell \ell^\top]$ .

We will treat  $\omega^2 \ell \ell^\top \in \mathbb{R}^{m \times m}$  as a  $d \times d$  matrix, because all but the top-left  $d \times d$  entries are zero. The expectation of the  $(l_1, l_2)$ -entry is given by

$$\begin{aligned} \mathbb{E}[\omega^2 \ell_{l_1} \ell_{l_2}] &= \mathbb{E}\left[\left(1 - \frac{\|\mathbf{x}\|^2}{s}\right) \ell_{l_1} \ell_{l_2}\right] \\ &= \mathbb{E}[\ell_{l_1} \ell_{l_2}] - \frac{1}{s} \mathbb{E}\left[\|\mathbf{x}\|^2 \ell_{l_1} \ell_{l_2}\right] \\ &= \mathbb{E}[\ell_{l_1} \ell_{l_2}] - \frac{1}{s} \mathbb{E}\left[\left(\|\ell\|^2 + \|\mathbf{c}\|^2 + \|\mathbf{e}\|^2 + 2 \langle \ell + \mathbf{c}, \mathbf{e} \rangle\right) \ell_{l_1} \ell_{l_2}\right]. \end{aligned}$$

Since the vector formed from the first  $d$  entries of  $\ell$  is uniformly generated in  $B_0^d(r)$ , we can write the  $j$ th coordinate of  $\ell$ ,  $1 \leq j \leq d$ , as

$$(A.11) \quad \ell_j = r u^{1/d} s_j,$$

where  $u \sim \text{unif}(0, 1)$  and  $s_j$  is the  $j$ th coordinate of  $\mathbf{s} \in \mathbb{R}^d$ , a random vector uniformly sampled from the unit sphere  $S^{d-1}$  [56]. It follows that

$$\mathbb{E}[\ell_{l_1} \ell_{l_2}] = \begin{cases} \frac{r^2}{d} \mathbb{E}[u^{2/d}] = \frac{r^2}{d+2}, & l_1 = l_2 \\ 0, & l_1 \neq l_2. \end{cases}$$

It remains to compute  $\mathbb{E}[\|\ell\|^2 \ell_{l_1} \ell_{l_2}]$ ,  $\mathbb{E}[\|\mathbf{c}\|^2 \ell_{l_1} \ell_{l_2}]$ ,  $\mathbb{E}[\|\mathbf{e}\|^2 \ell_{l_1} \ell_{l_2}]$ , and  $\mathbb{E}[\langle \ell + \mathbf{c}, \mathbf{e} \rangle \ell_{l_1} \ell_{l_2}]$ . Using the expansion in Eq. (A.11), we have

$$\mathbb{E}[\|\ell\|^2 \ell_{l_1} \ell_{l_2}] = \frac{r^4 d}{d+4} \sum_{j=1}^d \mathbb{E}[s_{l_1} s_{l_2} s_j^2]$$

and

$$\begin{aligned} \mathbb{E}[\|\mathbf{c}\|^2 \ell_{l_1} \ell_{l_2}] &= \mathbb{E}\left[\left(\frac{1}{4} \sum_{i=d+1}^m \left(\sum_{j_1=1}^d \kappa_{j_1}^{(i)} \ell_{j_1}^2\right)^2\right) \ell_{l_1} \ell_{l_2}\right] \\ &= \mathbb{E}\left[\left(\frac{r^4}{4} u^{4/d} \sum_{i=d+1}^m \left(\sum_{j_1=1}^d \kappa_{j_1}^{(i)} s_{j_1}^2\right)^2\right) r^2 u^{2/d} s_{l_1} s_{l_2}\right] \end{aligned}$$

## A.5. BOUNDING TERMS

---

$$\begin{aligned}
&= \frac{r^6 d}{4(d+6)} \mathbb{E} \left[ \left( \sum_{i=d+1}^m \left( \sum_{j_1=1}^d \kappa_{j_1}^{(i)} s_{j_1}^2 \right)^2 \right) s_{l_1} s_{l_2} \right] \\
&= \frac{r^6 d}{4(d+6)} \left[ \sum_{i=d+1}^m \left( \sum_{j_1=1}^d (\kappa_{j_1}^{(i)})^2 \mathbb{E}[s_{j_1}^4 s_{l_1} s_{l_2}] + \sum_{\substack{j_1, j_2=1 \\ j_1 \neq j_2}}^d \kappa_{j_1}^{(i)} \kappa_{j_2}^{(i)} \mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{l_1} s_{l_2}] \right) \right].
\end{aligned}$$

Since the distribution of  $\mathbf{s} = [s_1, \dots, s_d]^T$  is invariant to sign changes in each coordinate, the expectation of any product of coordinates of  $\mathbf{s}$  involving an odd power is 0. For example, for  $l_1 \neq l_2$ , we have that  $\mathbb{E}[s_{l_1}^5 s_{l_2}] = \mathbb{E}[s_{l_1}^5 (-s_{l_2})] = -\mathbb{E}[s_{l_1}^5 s_{l_2}] \Rightarrow \mathbb{E}[s_{l_1}^5 s_{l_2}] = 0$ . In the same manner, we obtain  $\mathbb{E}[s_{j_1}^4 s_{l_1} s_{l_2}] = \mathbb{E}[s_{j_1}^2 s_{l_1} s_{l_2}] = \mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{l_1} s_{l_2}] = 0$  for  $1 \leq l_1 \neq l_2 \leq d$ . Thus these terms reduce to

$$\mathbb{E}[\|\boldsymbol{\ell}\|^2 \ell_{l_1}^2] = \frac{r^4 d}{d+4} \sum_{j_1=1}^d \mathbb{E}[s_{j_1}^2 s_{l_1}^2]$$

for arbitrary  $1 \leq j_1 \leq d$ , and

$$\mathbb{E}[\|\boldsymbol{c}\|^2 \ell_{l_1}^2] = \frac{r^6 d}{4(d+6)} \left[ \sum_{i=d+1}^m \left( \sum_{j_1=1}^d (\kappa_{j_1}^{(i)})^2 \mathbb{E}[s_{j_1}^4 s_{l_1}^2] + \sum_{\substack{j_1, j_2=1 \\ j_1 \neq j_2}}^d \kappa_{j_1}^{(i)} \kappa_{j_2}^{(i)} \mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{l_1}^2] \right) \right]$$

for arbitrary  $1 \leq j_1 \neq j_2 \leq d$ .

For the remaining two terms in the expansion of  $\mathbb{E}[\omega^2 \ell_{l_1} \ell_{l_2}]$ , we have

$$\mathbb{E}[\|\boldsymbol{e}\|^2 \ell_{l_1} \ell_{l_2}] = \mathbb{E}[\|\boldsymbol{e}\|^2] \mathbb{E}[\ell_{l_1} \ell_{l_2}] = \frac{r^2 \sigma^2 m}{d+2}$$

for  $l_1 = l_2$  by the independence of  $\boldsymbol{\ell}$  and  $\boldsymbol{e}$  (and equal to zero otherwise), and

$$\mathbb{E}[\langle \boldsymbol{\ell} + \boldsymbol{c}, \boldsymbol{e} \rangle \ell_{l_1} \ell_{l_2}] = 0,$$

by the independence of  $\boldsymbol{\ell} + \boldsymbol{c}$  and  $\boldsymbol{e}$ .

Putting everything together, we conclude that  $\mathbb{E}[\omega^2 \ell_{l_1} \ell_{l_2}] = 0$  if  $l_1 \neq l_2$  and that

$$\begin{aligned}
\mathbb{E}[\omega^2 \ell_{l_1}^2] &= \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4 d}{d+4} \sum_{j_1=1}^d \mathbb{E}[s_{l_1}^2 s_{j_1}^2] \right. \\
&\quad \left. + \frac{r^6 d}{4(d+6)} \left[ \sum_{i=d+1}^m \left( \sum_{j_1=1}^d (\kappa_{j_1}^{(i)})^2 \mathbb{E}[s_{j_1}^4 s_{l_1}^2] + \sum_{\substack{j_1, j_2=1 \\ j_1 \neq j_2}}^d \kappa_{j_1}^{(i)} \kappa_{j_2}^{(i)} \mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{l_1}^2] \right) \right] \right)
\end{aligned}$$

## A.5. BOUNDING TERMS

---

$$\begin{aligned}
&= \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4 d}{d+4} \left( (d-1) \mathbb{E}[s_{l_1}^2 s_{j_1}^2] + \mathbb{E}[s_{l_1}^4] \right) \right. \\
&\quad \left. + \frac{r^6 d}{4(d+6)} \left[ \left( (d-1) \mathbb{E}[s_{j_1}^4 s_{l_1}^2] + \mathbb{E}[s_{l_1}^6] \right) (K^{(+)})^2 \right. \right. \\
&\quad \left. \left. + \left( 2(d-1) \mathbb{E}[s_{l_1}^4 s_{j_1}^2] + (d^2 - 3d + 2) \mathbb{E}[s_{l_1}^2 s_{j_1}^2 s_{j_2}^2] \right) \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right),
\end{aligned}$$

where  $j_1$  and  $j_2$  are arbitrary in  $1, \dots, d$  and  $j_1 \neq j_2 \neq l_1$ . In the last term, we utilized Kaslovsky and Meyer's notation [57]

$$(A.12) \quad K_{j_1 j_2}^{i_1 i_2} := \sum_{\substack{j_1, j_2=1 \\ j_1 \neq j_2}}^d \kappa_{j_1}^{(i_1)} \kappa_{j_2}^{(i_2)}$$

(note that  $K_{j_1 j_2}^{i_1 i_2}$  does not actually depend on  $j_1$  and  $j_2$ ). Also recall the definition of  $K^{(+)}$  given in Eq. (A.8).

It follows that we can set

$$\begin{aligned}
\lambda_{\min} \left( \sum_{k=1}^n \mathbb{E}[X_k] \right) &= \lambda_{\max} \left( \sum_{k=1}^n \mathbb{E}[X_k] \right) \\
&= \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4 d}{d+4} \left( (d-1) \mathbb{E}[s_{l_1}^2 s_{j_1}^2] + \mathbb{E}[s_{l_1}^4] \right) \right. \\
&\quad \left. + \frac{r^6 d}{4(d+6)} \left[ \left( (d-1) \mathbb{E}[s_{j_1}^4 s_{l_1}^2] + \mathbb{E}[s_{l_1}^6] \right) (K^{(+)})^2 \right. \right. \\
&\quad \left. \left. + \left( 2(d-1) \mathbb{E}[s_{l_1}^4 s_{j_1}^2] + (d^2 - 3d + 2) \mathbb{E}[s_{l_1}^2 s_{j_1}^2 s_{j_2}^2] \right) \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\
&= \frac{r^2}{d+2} \\
&\quad - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\
&=: \mu,
\end{aligned}$$

where  $\mu = \mu_{\min} = \mu_{\max}$ .

Let us set  $\delta := \xi_{\lambda_d^W} \sqrt{\frac{2\lambda_\infty}{\mu}}$ , for  $\xi_{\lambda_d^W}$  a probability constant that relates to our bound on  $\lambda_d^W$ , the smallest (nonzero) eigenvalue of  $(1/n)LD_0(LD_0)^\top$ . By Theorem A.5.1,  $\lambda_d^W$  satisfies  $\lambda_d^W > (1 - \delta)\mu$

## A.5. BOUNDING TERMS

---

with probability greater than  $1 - d \left[ \frac{e^{-\delta}}{(1-\delta)^{(1-\delta)}} \right]^{\frac{\mu}{\lambda_\infty}}$ . Expanding terms renders

$$\lambda_d^W > (1-\delta)\mu = \left( 1 - \xi_{\lambda_d^W} \frac{\sqrt{2\lambda_\infty}}{\sqrt{\mu}} \right) \mu = \mu - \xi_{\lambda_d^W} \sqrt{2\lambda_\infty} \sqrt{\mu},$$

and plugging in  $\mu$  and  $\lambda_\infty$  produces

$$\begin{aligned} \lambda_d^W &> \left[ \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2\sigma^2m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \right] \\ &\quad - \xi_{\lambda_d^W} \sqrt{2\lambda_\infty} \sqrt{\frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2\sigma^2m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)} \\ &= \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2\sigma^2m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\ &\quad - \xi_{\lambda_d^W} \frac{r^3\sqrt{2}}{\sqrt{n}} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)}. \end{aligned}$$

Next, we show that this bound on  $\lambda_d^W$  holds with probability greater than  $1 - de^{\xi_{\lambda_d^W}^2}$  by showing that

$$1 - d \left[ \frac{e^{-\delta}}{(1-\delta)^{(1-\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} \geq 1 - de^{\xi_{\lambda_d^W}^2}.$$

We have that

$$\begin{aligned} 1 - d \left[ \frac{e^{-\delta}}{(1-\delta)^{(1-\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} &\geq 1 - de^{-\xi_{\lambda_d^W}^2} \\ \Leftrightarrow \left[ \frac{e^{-\delta}}{(1-\delta)^{(1-\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} &\leq e^{-\xi_{\lambda_d^W}^2} \\ \Leftrightarrow \left[ \frac{e^{-\delta}}{(1-\delta)^{(1-\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} &\leq e^{-\frac{\delta^2}{2} \frac{\mu}{\lambda_\infty}} \\ \Leftrightarrow \frac{e^{-\delta}}{(1-\delta)^{(1-\delta)}} &\leq e^{-\frac{\delta^2}{2}}, \end{aligned}$$

which is true for  $\delta \in [0, 1]$ . Note that this requires that

$$(A.13) \quad \xi_{\lambda_d^W} \leq \frac{\sqrt{n}}{\sqrt{2}} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2 m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)}.$$

We state this result formally:

LEMMA A.5.1. *For probability constant  $\xi_{\lambda_d^W}$  obeying Eq. (A.13), the smallest nonzero eigenvalue of  $(1/n)LD_0(LD_0)^T$  satisfies*

$$\begin{aligned} \lambda_d^W &> \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\ &\quad - \xi_{\lambda_d^W} \frac{r^3 \sqrt{2}}{\sqrt{n}} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2 m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)} \end{aligned}$$

with probability greater than  $1 - de^{-\xi_{\lambda_d^W}^2}$ , where  $K^{(+)}$  is given by Eq. (A.8) and  $K_{j_1 j_2}^{ii}$  is given by Eq. (A.12).

Now we will produce an upper bound for the largest eigenvalue,  $\lambda_1^W$ , of  $(1/n)LD_0(LD_0)^T$ . Set  $\delta := \xi_{\lambda_1^W}(5/2)\sqrt{\lambda_\infty/\mu}$ . Theorem A.5.1 states that  $\lambda_1^W$  satisfies  $\lambda_1^W < (1 + \delta)\mu$  with probability at least  $1 - d \left[ \frac{e^{-\delta}}{(1+\delta)^{(1+\delta)}} \right]^{\frac{\mu}{\lambda_\infty}}$ . Expanding terms gives

$$\begin{aligned} \lambda_1^W &< (1 + \delta)\mu \\ &= \left( 1 + \xi_{\lambda_1^W} \frac{5}{2} \frac{\sqrt{\lambda_\infty}}{\sqrt{\mu}} \right) \mu \\ &= \mu + \xi_{\lambda_1^W} \frac{5}{2} \sqrt{\lambda_\infty} \sqrt{\mu} \\ &= \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\ &\quad + \xi_{\lambda_1^W} \frac{5r^3}{2\sqrt{n}} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2 m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)} \end{aligned}$$

## A.5. BOUNDING TERMS

---

Lastly, this bound on  $\lambda_1^W$  holds with probability greater than  $1 - de^{-\xi_{\lambda_1^W}^2}$  since

$$\begin{aligned} 1 - d \left[ \frac{e^{-\delta}}{(1 + \delta)^{(1+\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} &> 1 - de^{-\xi_{\lambda_1^W}^2} \\ \Leftrightarrow \left[ \frac{e^{-\delta}}{(1 + \delta)^{(1+\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} &< e^{-\xi_{\lambda_1^W}^2} \\ \Leftrightarrow \left[ \frac{e^{-\delta}}{(1 + \delta)^{(1+\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} &< e^{-\delta^2 \frac{4}{25} \frac{\mu}{\lambda_\infty}} \\ \Leftrightarrow \left[ \frac{e^{-\delta}}{(1 + \delta)^{(1+\delta)}} \right]^{\frac{\mu}{\lambda_\infty}} &< e^{-\frac{4}{25} \delta^2} \end{aligned}$$

is true for all  $\delta, \xi_{\lambda_1} \geq 0$ .

LEMMA A.5.2. *The largest eigenvalue of  $(1/n)LD_0(LD_0)^\top$  satisfies*

$$\begin{aligned} \lambda_1^W &< \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\ &\quad + \xi_{\lambda_1^W} \frac{5r^3}{2\sqrt{n}} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2 m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)} \end{aligned}$$

with probability greater than  $1 - de^{-\xi_{\lambda_1^W}^2}$ , where  $K^{(+)}$  is given by Eq. (A.8) and  $K_{j_1 j_2}^{ii}$  is given by Eq. (A.12).

We finish this section by setting

$$\begin{aligned} (\text{A.14}) \quad \lambda_{\text{bound}}^W(\xi) := & \frac{r^2}{d+2} - \frac{1}{s} \left( \frac{r^2 \sigma^2 m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\ & + \xi \frac{5r^3}{2\sqrt{n}} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2 m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)} \end{aligned}$$

for later use. By design, all the eigenvalues of  $(1/n)LD_0(LD_0)^\top$  are no greater than  $\lambda_{\text{bound}}^W(\xi)$  with probability depending on  $\xi$ .

**A.5.2. Curvature Eigenvalue Bounds.** We postpone this analysis until the end of Section A.5.4.

## A.5. BOUNDING TERMS

---

**A.5.3. Noise Eigenvalue Bounds.** To bound the eigenvalues of  $W_i^\top (1/n) E D_0 (E D_0)^\top W_i$  for  $i = 1, 2$ , we will use the following theorem:

**THEOREM A.5.2** (Edelman [35], Vershynin [97]). *Let  $A$  be an  $m \times n$  matrix,  $m \leq n$ , whose entries are independent standard normal random variables. Then for every  $t \geq 0$ , with probability at least  $1 - 2\exp(-t^2/2)$ , we have*

$$\sqrt{n} - \sqrt{m} - t \leq s_{\min}(A) \leq s_{\max}(A) \leq \sqrt{n} + \sqrt{m} + t,$$

where  $s_{\min}(A)$  and  $s_{\max}(A)$  denote the smallest and largest singular values of  $A$ , respectively.

In order to apply Theorem A.5.2, we first determine a constant  $\alpha$  such that  $\alpha E D_0$  has standard normal entries. This is a critical step in bounding all the terms involving the noise matrix  $E$ , as we will see.

Recall that  $E D_0 = [d_1 e^{(1)}, \dots, d_n e^{(n)}] \in \mathbb{R}^{m \times n}$ , where the entries of the vectors  $e^{(k)}$  have Gaussian distribution with mean 0 and variance  $\sigma^2$ . For  $\mathbf{x}$  the random variable from which the samples  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$  are realizations, let  $\mathbf{e}$  denote its random noise component. As before, we define  $\omega$  to be the random variable denoting the weight  $d_k$  from Eq. (A.10). Thus  $\mathbf{x}^{(k)}$  and  $d_k e^{(k)}$  are the  $k$ th realizations of the random vectors  $\mathbf{x}$  and  $\omega \mathbf{e}$ , respectively. Since an arbitrary coordinate  $e_l$  of  $\mathbf{e}$  satisfies  $e_l \sim \mathcal{N}(0, \sigma^2)$ , the scaled noise component  $\omega \mathbf{e}$  also has Gaussian-distributed entries with mean 0. Additionally, since  $\omega$  and  $e_l$  are independent, we have that

$$\begin{aligned} \text{var}(\alpha(\omega e_l)) &= \alpha^2 \text{var}(\omega e_l) = \alpha^2 (\mathbb{E}[(\omega e_l)^2] - \mathbb{E}[\omega e_l]^2) \\ &= \alpha^2 (\mathbb{E}[\omega^2] \mathbb{E}[e_l^2] - \mathbb{E}[\omega]^2 \mathbb{E}[e_l]^2) \\ &= \alpha^2 \sigma^2 \mathbb{E}[\omega^2]. \end{aligned}$$

We compute

$$\begin{aligned} \mathbb{E}[\omega^2] &= 1 - \frac{1}{s} \mathbb{E}[\|\mathbf{x}\|^2] = 1 - \frac{1}{s} \left( \mathbb{E}[\|\ell\|^2] + \mathbb{E}[\|\mathbf{c}\|^2] + \mathbb{E}[\|e\|^2] + 2\mathbb{E}[\langle \ell, \mathbf{c} \rangle] + 2\mathbb{E}[\langle \ell + \mathbf{c}, e \rangle] \right) \\ &= 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right), \end{aligned}$$

## A.5. BOUNDING TERMS

---

where

$$(A.15) \quad \gamma := \frac{\sum_{i=d+1}^m (3K_{j_1 j_1}^{ii} + K_{j_1 j_2}^{ii})}{2(d+2)(d+4)}.$$

This is a quantity defined in Kaslovsky and Meyer's original analysis [57]. It follows that setting

$$(A.16) \quad \alpha := \left( \sigma^2 \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \right)^{-1/2}$$

will result in  $\alpha E D_0$  having standard normal entries.

In order to bound the noise eigenvalues, we view the columns of  $E$  as  $n$  realizations of the random noise vector  $\mathbf{e} \in \mathbb{R}^m$  and partition  $\mathbf{e} = [\mathbf{e}_1^\top, \mathbf{e}_2^\top]^\top$  for  $\mathbf{e}_1 \in \mathbb{R}^d$  and  $\mathbf{e}_2 \in \mathbb{R}^{m-d}$ . Then  $W_i^\top (1/n) E D_0 (E D_0)^\top W_i$  depends only on the realizations of  $\mathbf{e}_1$  for  $i = 1$  and those of  $\mathbf{e}_2$  for  $i = 2$ . Next, define  $Z := \alpha E D_0$  for  $\alpha$  given by Eq. (A.16). By Theorem A.5.2 and standard inequalities, we have

$$\begin{aligned} \lambda_{\max} \left( \frac{1}{n} W_1^\top Z Z^\top W_1 \right) &\leq \left( 1 - \frac{1}{n} \right)^{-1} \left( 1 + \frac{5}{2\sqrt{n}} (\sqrt{d} + \xi_{\mathbf{e}_1} \sqrt{2}) \right) \\ \Rightarrow \lambda_{\max} \left( \frac{1}{n} W_1^\top E D_0 (E D_0)^\top W_1 \right) &\leq \frac{1}{\alpha^2} \left( 1 - \frac{1}{n} \right)^{-1} \left( 1 + \frac{5}{2\sqrt{n}} (\sqrt{d} + \xi_{\mathbf{e}_1} \sqrt{2}) \right) \\ \Rightarrow \lambda_{\max} \left( \frac{1}{n} W_1^\top E D_0 (E D_0)^\top W_1 \right) &\leq \sigma^2 \left( 1 - \frac{1}{s} \left[ \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right] \right) \left( 1 - \frac{1}{n} \right)^{-1} \left( 1 + \frac{5}{2\sqrt{n}} (\sqrt{d} + \xi_{\mathbf{e}_1} \sqrt{2}) \right) \end{aligned}$$

with probability at least  $1 - e^{-\xi_{\mathbf{e}_1}^2}$  for  $n > 4(\sqrt{d} + \xi_{\mathbf{e}_1})$ , and similarly

$$\lambda_{\max} \left( \frac{1}{n} W_2^\top E D_0 (E D_0)^\top W_2 \right) \leq \sigma^2 \left( 1 - \frac{1}{s} \left[ \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right] \right) \left( 1 - \frac{1}{n} \right)^{-1} \left( 1 + \frac{5}{2\sqrt{n}} (\sqrt{m-d} + \xi_{\mathbf{e}_2} \sqrt{2}) \right)$$

with probability at least  $1 - e^{-\xi_{\mathbf{e}_2}^2}$  for  $n > 4(\sqrt{m-d} + \xi_{\mathbf{e}_2})$ . We state these results formally:

**LEMMA A.5.3.** *The largest eigenvalue of  $W_1^\top (1/n) E D_0 (E D_0)^\top W_1$  satisfies*

$$\lambda_{\max} \left( \frac{1}{n} W_1^\top E D_0 (E D_0)^\top W_1 \right) \leq \sigma^2 \left( 1 - \frac{1}{s} \left[ \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right] \right) \left( 1 - \frac{1}{n} \right)^{-1} \left( 1 + \frac{5}{2\sqrt{n}} (\sqrt{d} + \xi_{\mathbf{e}_1} \sqrt{2}) \right)$$

*with probability at least  $1 - 2e^{-\xi_{\mathbf{e}_1}^2}$  for  $n > 4(\sqrt{d} + \xi_{\mathbf{e}_1})$ . The quantity  $\gamma$  is defined in Eq. (A.15), and the notation “ $\xi_{\mathbf{e}_1}$ ” indicates that  $W_1^\top E D_0 (E D_0)^\top W_1$  only depends on the first  $d$  coordinates of the realizations of the random noise vector  $\mathbf{e} \in \mathbb{R}^m$ .*

## A.5. BOUNDING TERMS

---

LEMMA A.5.4. *The largest eigenvalue of  $W_2^\top (1/n)ED_0(ED_0)^\top W_2$  satisfies*

$$\lambda_{\max}\left(\frac{1}{n}W_2^\top ED_0(ED_0)^\top W_2\right) \leq \sigma^2 \left(1 - \frac{1}{s} \left[\frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m\right]\right) \left(1 - \frac{1}{n}\right)^{-1} \left(1 + \frac{5}{2\sqrt{n}} (\sqrt{m-d} + \xi_{\mathbf{e}_2} \sqrt{2})\right)$$

*with probability at least  $1 - 2e^{-\xi_{\mathbf{e}_2}^2}$  for  $n > 4(\sqrt{m-d} + \xi_{\mathbf{e}_2})$ . The quantity  $\gamma$  is defined in Eq. (A.15), and the notation “ $\xi_{\mathbf{e}_2}$ ” indicates that  $W_2^\top ED_0(ED_0)^\top W_2$  only depends on the last  $m-d$  coordinates of the realizations of the random noise vector  $\mathbf{e} \in \mathbb{R}^m$ .*

**A.5.4. Curvature Bounds.** Our goal in this section is to bound

$$\|W_2^\top (1/n)CD_0(CD_0)^\top W_2\|_F.$$

Since  $W_2$  contains the last  $m-d$  eigenvectors of  $(1/n)LD_0(LD_0)^\top$ , we have that

$$W_2 = \begin{bmatrix} 0 & 0 \\ 0 & I_{(m-d) \times (m-d)} \end{bmatrix} \in \mathbb{R}^{m \times (m-d)},$$

because  $LD_0(LD_0)^\top$  contains nonzeros only in its first  $d$  rows and columns. Further, since  $CD_0(CD_0)^\top$  contains nonzeros only in its *last*  $m-d$  rows and columns, it follows that  $W_2^\top CD_0(CD_0)^\top W_2$  simply extracts the nonzero block of  $CD_0(CD_0)^\top$ . Consequently,  $CD_0(CD_0)^\top$  and  $W_2^\top CD_0(CD_0)^\top W_2$  have the same Frobenius norm.

We can view the quantity  $(1/n)CD_0(CD_0)^\top$  as the empirical expectation of  $n$  samples of the scaled random curvature vector  $\omega \mathbf{c}$ , where  $\omega$  again denotes the random weight variable given by Eq. (A.10). We write

$$\frac{1}{n}CD_0(CD_0)^\top = \hat{\mathbb{E}}[\omega^2 \mathbf{c} \mathbf{c}^\top],$$

and we use the following theorem:

**THEOREM A.5.3** (Shawe-Taylor and Cristianini [84]). *Given  $n$  realizations of a random matrix  $Y$  distributed with probability  $\mathbb{P}_Y$ , we have*

$$\mathbb{P}_Y \left\{ \|\mathbb{E}[Y] - \hat{\mathbb{E}}[Y]\|_F \leq \frac{R}{\sqrt{n}} (2 + \xi \sqrt{2}) \right\} \geq 1 - e^{-\xi^2},$$

where  $R := \sup_{\text{supp}(\mathbb{P}_Y)} \|Y\|_F$ .

## A.5. BOUNDING TERMS

---

Since

$$\left| \|\mathbb{E}[W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2]\|_F - \|\hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2]\|_F \right| \leq \left\| \mathbb{E}[W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2] - \hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2] \right\|_F,$$

setting  $Y := W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2$  in Theorem A.5.3 renders

$$\left| \|\mathbb{E}[W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2]\|_F - \|\hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2]\|_F \right| \leq \frac{R_c}{\sqrt{n}} (2 + \xi\sqrt{2})$$

with probability at least  $1 - e^{-\xi^2}$ , where  $R_c$  is an upper bound on  $\|W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2\|_F$  (note that it is unnecessary to use an exact supremum in Theorem A.5.3). Because  $\|W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2\|_F \leq \|W_2^\top \omega \mathbf{c}\|_F \|\omega \mathbf{c}^\top W_2\|_F = \|W_2^\top \omega \mathbf{c}\|_F^2$ , it suffices to bound this last quantity. Observe

$$\begin{aligned} \|W_2^\top \omega \mathbf{c}\|_F^2 &= \|\omega \mathbf{c}\|^2 \leq \|\mathbf{c}\|^2 = \frac{1}{4} \sum_{i=d+1}^m \left( \kappa_1^{(i)} \ell_1^2 + \dots + \kappa_d^{(i)} \ell_d^2 \right)^2 \leq \frac{r^4}{4} \sum_{i=d+1}^m \left( \sum_{j=1}^d (\kappa_j^{(i)})^2 \right)^2 \\ &= \frac{r^4}{4} \sum_{i=d+1}^m (K_i^{(+)})^2 \\ &\leq \frac{(K^{(+)})^2 r^4}{4}, \end{aligned}$$

and so we set  $R_c := \frac{K^{(+)} r^2}{2}$ . The first inequality follows by Cauchy-Schwarz, and the terms  $K_i^{(+)}$  and  $K^{(+)}$  are the curvature constants defined in Eq. (A.8).

We next need to compute the true expectation  $\mathbb{E}[W_2^\top \omega^2 \mathbf{c} \mathbf{c}^\top W_2]$ , which requires computing  $\mathbb{E}[\omega^2 c_{l_1} c_{l_2}]$  for arbitrary  $d+1 \leq l_1, l_2 \leq m$ . We expand

$$\begin{aligned} \mathbb{E}[\omega^2 c_{l_1} c_{l_2}] &= \mathbb{E}\left[\left(1 - \frac{\|\mathbf{x}\|^2}{s}\right) c_{l_1} c_{l_2}\right] \\ &= \mathbb{E}[c_{l_1} c_{l_2}] - \mathbb{E}\left[\frac{\|\mathbf{x}\|^2}{s} c_{l_1} c_{l_2}\right]. \end{aligned}$$

From Kaslovsky and Meyer's original analysis [57], we have that

$$\mathbb{E}[c_{l_1} c_{l_2}] = \frac{[3K_{j_1 j_1}^{l_1 l_2} + K_{j_1 j_2}^{l_1 l_2}]r^4}{4(d+2)(d+4)},$$

## A.5. BOUNDING TERMS

---

where  $K_{j_1 j_2}^{l_1 l_2}$  is given by Eq. (A.12) and

$$(A.17) \quad K_{j_1 j_1}^{l_1 l_2} := \sum_{j=1}^d \kappa_{j_1}^{(l_1)} \kappa_{j_1}^{(l_2)}.$$

Expanding the second term, we have

$$\begin{aligned} \mathbb{E}\left[\frac{\|\boldsymbol{x}\|^2}{s} c_{l_1} c_{l_2}\right] &= \frac{1}{s} \mathbb{E}\left[\left(\|\boldsymbol{\ell}\|^2 + \|\boldsymbol{c}\|^2 + \|\boldsymbol{e}\|^2 + 2\langle \boldsymbol{\ell} + \boldsymbol{c}, \boldsymbol{e} \rangle\right) c_{l_1} c_{l_2}\right] \\ &= \frac{1}{s} \left( \mathbb{E}[\|\boldsymbol{\ell}\|^2 c_{l_1} c_{l_2}] + \mathbb{E}[\|\boldsymbol{c}\|^2 c_{l_1} c_{l_2}] + \mathbb{E}[\|\boldsymbol{e}\|^2 c_{l_1} c_{l_2}] + 2\mathbb{E}[\langle \boldsymbol{\ell} + \boldsymbol{c}, \boldsymbol{e} \rangle c_{l_1} c_{l_2}] \right). \end{aligned}$$

We compute these quantities one at a time. For the first, we have

$$\begin{aligned} \mathbb{E}[\|\boldsymbol{\ell}\|^2 c_{l_1} c_{j_1}] &= \mathbb{E}\left[\left(\sum_{j_1=1}^d \ell_{j_1}^2\right)\left(\sum_{j_2=1}^d \kappa_{j_2}^{(l_1)} \kappa_{j_2}^{(l_2)} \ell_{j_2}^4 + \sum_{\substack{j_2, j_3=1 \\ j_2 \neq j_3}}^d \kappa_{j_2}^{(l_1)} \kappa_{j_3}^{(l_2)} \ell_{j_2}^2 \ell_{j_3}^2\right)\right] \\ &= \mathbb{E}\left[\sum_{j_1=1}^d \sum_{j_2=1}^d \kappa_{j_2}^{(l_1)} \kappa_{j_2}^{(l_2)} \ell_{j_1}^2 \ell_{j_2}^4\right] + \mathbb{E}\left[\sum_{j_1=1}^d \sum_{\substack{j_2, j_3=1 \\ j_2 \neq j_3}}^d \kappa_{j_2}^{(l_1)} \kappa_{j_3}^{(l_2)} \ell_{j_1}^2 \ell_{j_2}^2 \ell_{j_3}^2\right] \\ &= \frac{dr^6}{d+6} \left( \sum_{j_1=1}^d \sum_{j_2=1}^d \kappa_{j_2}^{(l_1)} \kappa_{j_2}^{(l_2)} \mathbb{E}[s_{j_1}^2 s_{j_2}^4] + \sum_{j_1=1}^d \sum_{\substack{j_2, j_3=1 \\ j_2 \neq j_3}}^d \kappa_{j_2}^{(l_1)} \kappa_{j_3}^{(l_2)} \mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{j_3}^2] \right) \\ &= \frac{dr^6}{d+6} \left( \left[ d\mathbb{E}[s_{j_1}^6] + (d^2 - d)\mathbb{E}[s_{j_1}^2 s_{j_2}^4] \right] K_{j_1 j_1}^{l_1 l_2} \right. \\ &\quad \left. + \left[ 2d(d-1)\mathbb{E}[s_{j_1}^4 s_{j_2}^2] + (d-2)(d^2 - d)\mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{j_3}^2] \right] K_{j_2 j_3}^{l_1 l_2} \right), \end{aligned}$$

for  $1 \leq j_1 \neq j_2 \neq j_3 \leq d$ . Plugging in the expectations, this reduces to

$$\mathbb{E}[\|\boldsymbol{\ell}\|^2 c_{l_1} c_{j_1}] = \frac{dr^6 K_{j_1 j_1}^{l_1 l_2}}{d+6}.$$

For the curvature term, we have

$$\mathbb{E}[\|\boldsymbol{c}\|^2 c_{l_1} c_{l_2}] = \frac{1}{16} \mathbb{E}\left[\left(\sum_{i=d+1}^m \left[\sum_{j_1=1}^d \kappa_{j_1}^{(i)} \ell_{j_1}^2\right]^2\right) \left(\sum_{j_2=1}^d \kappa_{j_2}^{(l_1)} \kappa_{j_2}^{(l_2)} \ell_{j_2}^4 + \sum_{\substack{j_2, j_3=1 \\ j_2 \neq j_3}}^d \kappa_{j_2}^{(l_1)} \kappa_{j_3}^{(l_2)} \ell_{j_2}^2 \ell_{j_3}^2\right)\right]$$

## A.5. BOUNDING TERMS

---

$$\begin{aligned}
&= \frac{1}{16} \mathbb{E} \left[ \left( \sum_{i=d+1}^m \sum_{j_1=1}^d (\kappa_{j_1}^{(i)})^2 \ell_{j_1}^4 + \sum_{i=d+1}^m \sum_{\substack{j_4, j_5=1 \\ j_4 \neq j_5}}^d \kappa_{j_4}^{(i)} \kappa_{j_5}^{(i)} \ell_{j_4}^2 \ell_{j_5}^2 \right) \right. \\
&\quad \times \left. \left( \sum_{j_2=1}^d \kappa_{j_2}^{(l_1)} \kappa_{j_2}^{(l_2)} \ell_{j_2}^4 + \sum_{\substack{j_2, j_3=1 \\ j_2 \neq j_3}}^d \kappa_{j_2}^{(l_1)} \kappa_{j_3}^{(l_2)} \ell_{j_2}^2 \ell_{j_3}^2 \right) \right] \\
&= \frac{1}{16} \mathbb{E} \left[ \sum_{i=d+1}^m \sum_{j_1=1}^d (\kappa_{j_1}^{(i)})^2 \ell_{j_1}^4 \sum_{j_2=1}^d \kappa_{j_2}^{(l_1)} \kappa_{j_2}^{(l_2)} \ell_{j_2}^4 + \sum_{i=d+1}^m \sum_{\substack{j_4, j_5=1 \\ j_4 \neq j_5}}^d \kappa_{j_4}^{(i)} \kappa_{j_5}^{(i)} \ell_{j_4}^2 \ell_{j_5}^2 \sum_{j_2=1}^d \kappa_{j_2}^{(l_1)} \kappa_{j_2}^{(l_2)} \ell_{j_2}^4 \right. \\
&\quad \left. + \sum_{i=d+1}^m \sum_{j_1=1}^d (\kappa_{j_1}^{(i)})^2 \ell_{j_1}^4 \sum_{\substack{j_2, j_3=1 \\ j_2 \neq j_3}}^d \kappa_{j_2}^{(l_1)} \kappa_{j_3}^{(l_2)} \ell_{j_2}^2 \ell_{j_3}^2 + \sum_{i=d+1}^m \sum_{\substack{j_4, j_5=1 \\ j_4 \neq j_5}}^d \kappa_{j_4}^{(i)} \kappa_{j_5}^{(i)} \ell_{j_4}^2 \ell_{j_5}^2 \sum_{\substack{j_2, j_3=1 \\ j_2 \neq j_3}}^d \kappa_{j_2}^{(l_1)} \kappa_{j_3}^{(l_2)} \ell_{j_2}^2 \ell_{j_3}^2 \right],
\end{aligned}$$

which further reduces to

$$\begin{aligned}
\mathbb{E}[\|\mathbf{c}\|^2 c_{l_1} c_{l_2}] &= \frac{dr^8}{16(d+8)} \left( \left[ d \mathbb{E}[s_{j_1}^8] + (d^2 - d) \mathbb{E}[s_{j_1}^4 s_{j_2}^4] \right] (K^{(+)})^2 K_{j_1 j_1}^{l_1 l_2} \right. \\
&\quad + \left[ 2d(d-1) \mathbb{E}[s_{j_1}^6 s_{j_2}^2] + (d-2)(d^2 - d) \mathbb{E}[s_{j_1}^4 s_{j_2}^2 s_{j_3}^2] \right] K_{j_1 j_1}^{l_1 l_2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \\
&\quad + \left[ 2d(d-1) \mathbb{E}[s_{j_1}^6 s_{j_2}^2] + (d-2)(d^2 - d) \mathbb{E}[s_{j_1}^4 s_{j_2}^2 s_{j_3}^2] \right] (K^{(+)})^2 K_{j_1 j_2}^{l_1 l_2} \\
&\quad + \left[ 2d(d-1) \mathbb{E}[s_{j_1}^4 s_{j_2}^4] + (d^2 - d)(d-2)(d-3) \mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{j_3}^2 s_{j_4}^2] \right. \\
&\quad \left. + ((d^2 - d)^2 - 2d(d-1) - (d^2 - d)(d-2)(d-3)) \mathbb{E}[s_{j_1}^4 s_{j_2}^2 s_{j_3}^2] \right] K_{j_1 j_2}^{l_1 l_2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right),
\end{aligned}$$

where  $1 \leq j_1 \neq j_2 \neq j_3 \neq j_4 \leq d$  (recall that the subscript indexing for the curvature constants is arbitrary). Plugging in the expectations renders

$$\mathbb{E}[\|\mathbf{c}\|^2 c_{l_1} c_{l_2}] = \frac{dr^8 K_{j_1 j_1}^{l_1 l_2}}{16(d+2)(d+8)} \left( 3(K^{(+)})^2 + \frac{(d-1)(5d^2 - 10d + 48)}{(d+4)(d+6)} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right).$$

For the noise term, we have

$$\mathbb{E}[\|\mathbf{e}\|^2 c_{l_1} c_{l_2}] = \mathbb{E}[\|\mathbf{e}\|^2] \mathbb{E}[c_{l_1} c_{l_2}] = d\sigma^2 \frac{[3K_{j_1 j_1}^{l_1 l_2} + K_{j_1 j_2}^{l_1 l_2}] r^4}{4(d+2)(d+4)}$$

for  $1 \leq j_1 \neq j_2 \leq d$ , and finally, for the last term we have

$$\begin{aligned}
\mathbb{E}[\langle \boldsymbol{\ell} + \mathbf{c}, \mathbf{e} \rangle c_{l_1} c_{l_2}] &= \mathbb{E}\left[\left(\sum_{j=1}^d \ell_j e_j + \sum_{i=d+1}^m c_i e_i\right) c_{l_1} c_{l_2}\right] \\
&= \sum_{j=1}^d \mathbb{E}[\ell_j e_j c_{l_1} c_{l_2}] + \sum_{i=d+1}^m \mathbb{E}[c_i e_i c_{l_1} c_{l_2}] \\
&= \sum_{j=1}^d \mathbb{E}[e_j] \mathbb{E}[\ell_j c_{l_1} c_{l_2}] + \sum_{i=d+1}^m \mathbb{E}[e_i] \mathbb{E}[c_i c_{l_1} c_{l_2}] \\
&= 0.
\end{aligned}$$

Considering all the terms together, the true expectation satisfies

$$\begin{aligned}
\mathbb{E}[\omega^2 c_{l_1} c_{l_2}] &= \mathbb{E}[c_{l_1} c_{l_2}] - \frac{1}{s} \mathbb{E}[\|\boldsymbol{\ell}\|^2 c_{l_1} c_{l_2}] + \mathbb{E}[\|\mathbf{c}\|^2 c_{l_1} c_{l_2}] + \mathbb{E}[\|\mathbf{e}\|^2 c_{l_1} c_{l_2}] + 2\mathbb{E}[\langle \boldsymbol{\ell} + \mathbf{c}, \mathbf{e} \rangle c_{l_1} c_{l_2}] \\
&= \left(1 - \frac{d\sigma^2}{s}\right) \frac{[3K_{j_1 j_1}^{l_1 l_2} + K_{j_1 j_2}^{l_1 l_2}]r^4}{4(d+2)(d+4)} \\
&\quad - \frac{dr^6 K_{j_1 j_1}^{l_1 l_2}}{s} \left[ \frac{1}{d+6} + \frac{r^2}{16(d+2)(d+8)} \left( 3(K^{(+)})^2 + \frac{(d-1)(5d^2-10d+48)}{(d+4)(d+6)} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right) \right].
\end{aligned}$$

Theorem A.5.3 implies that

$$\left\| W_2^\top \left( \frac{1}{n} CD_0 (CD_0)^\top \right) W_2 \right\|_F \leq \left( \sum_{i,j=d+1}^m (\mathbb{E}[\omega^2 c_i c_j])^2 \right)^{1/2} + \frac{K^{(+)}(2 + \xi\sqrt{2})}{2\sqrt{n}}$$

with probability at least  $1 - e^{-\xi^2}$ . We set  $\xi = \xi_c$  to identify the probability constant corresponding to this curvature bound. Note that the above inequality also holds with probability *strictly* greater than  $1 - 2e^{-\xi_c^2}$ . Formally, we have

LEMMA A.5.5. *The Frobenius norm of the curvature term satisfies*

$$\left\| W_2^\top \left( \frac{1}{n} CD_0 (CD_0)^\top \right) W_2 \right\|_F \leq \left( \sum_{l_1, l_2=d+1}^m (\mathbb{E}[\omega^2 c_{l_1} c_{l_2}])^2 \right)^{1/2} + \frac{K^{(+)}(2 + \xi_c\sqrt{2})}{2\sqrt{n}}$$

## A.5. BOUNDING TERMS

---

with probability greater than  $1 - 2e^{-\xi_c^2}$ , where

$$(A.18) \quad \begin{aligned} \mathbb{E}[\omega^2 c_{l_1} c_{l_2}] &= \left(1 - \frac{d\sigma^2}{s}\right) \frac{[3K_{j_1 j_1}^{l_1 l_2} + K_{j_1 j_2}^{l_1 l_2}]r^4}{4(d+2)(d+4)} - \frac{dr^6 K_{j_1 j_1}^{l_1 l_2}}{s} \left[ \frac{1}{d+6} \right. \\ &\quad \left. + \frac{r^2}{16(d+2)(d+8)} \left( 3(K^{(+)})^2 + \frac{(d-1)(5d^2 - 10d + 48)}{(d+4)(d+6)} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right) \right]. \end{aligned}$$

Here,  $K^{(+)}$ ,  $K_{j_1 j_2}^{ii}$ , and  $K_{j_1 j_1}^{l_1 l_2}$  are defined in Eq. (A.8), Eq. (A.12), and Eq. (A.17), respectively.

To obtain a bound on  $\gamma_1^W$ , the largest eigenvalue of  $(1/n)CD_0(CD_0)^\top$ , we use that

$$\begin{aligned} \gamma_1^W &= \left\| \frac{1}{n} CD_0(CD_0)^\top \right\|_2 \\ &\leq \left\| \frac{1}{n} CD_0(CD_0)^\top \right\|_F \\ &= \left\| W_2^\top \left( \frac{1}{n} CD_0(CD_0)^\top \right) W_2 \right\|_F \\ &\leq \left( \sum_{i,j=d+1}^m (\mathbb{E}[\omega^2 c_i c_j])^2 \right)^{1/2} + \frac{K^{(+)}(2 + \xi\sqrt{2})}{2\sqrt{n}} \end{aligned}$$

with probability greater than  $1 - 2e^{-\xi^2}$ . This produces the following definition:

$$(A.19) \quad \gamma_{\text{bound}}^W(\xi) := \left( \sum_{i,j=d+1}^m (\mathbb{E}[\omega^2 c_i c_j])^2 \right)^{1/2} + \frac{K^{(+)}(2 + \xi\sqrt{2})}{2\sqrt{n}}.$$

**A.5.5. Curvature-Linear Bounds.** Next we bound  $\|W_2^\top (1/n)CD_0(LD_0)^\top W_1\|_F$  using the same approach as in the previous section. We observe that

$$\frac{1}{n} W_2^\top CD_0(LD_0)^\top W_1 = \hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c} \ell^\top W_1]$$

and use Theorem A.5.3 to obtain

$$\left| \|\mathbb{E}[W_2^\top \omega^2 \mathbf{c} \ell^\top W_1]\|_F - \|\hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c} \ell^\top W_1]\|_F \right| \leq \frac{R}{\sqrt{n}} (2 + \xi\sqrt{2})$$

with probability at least  $1 - e^{-\xi^2}$  and  $R \geq \sup \|W_2^\top \omega^2 \mathbf{c} \ell^\top W_1\|_F$ . Since

$$\|W_2^\top \omega^2 \mathbf{c} \ell^\top W_1\|_F \leq \|W_2^\top \omega \mathbf{c}\|_F \|\omega \ell^\top W_1\|_F \leq \|W_2^\top \mathbf{c}\|_F \|\ell^\top W_1\|_F,$$

we can use  $R_c$  from the previous section and set  $R_\ell := r$ . Thus we obtain  $R := R_c R_\ell$ .

## A.5. BOUNDING TERMS

---

To complete the bound, we will show that  $\|\mathbb{E}[W_2^\top \omega^2 \mathbf{c} \ell^\top W_1]\|_F = 0$ . Let us first consider  $\mathbb{E}[\omega^2 \mathbf{c} \ell^\top]$ . Since the matrix  $\omega^2 \mathbf{c} \ell^\top$  has nonzero entries  $\omega^2 c_{l_1} \ell_{l_2}$  for  $l_1 = d+1, \dots, m$  and  $l_2 = 1, \dots, d$ , we have that

$$\begin{aligned}\omega^2 c_{l_1} \ell_{l_2} &= \frac{1}{2} \omega^2 \left( \kappa_1^{(l_1)} \ell_1^2 + \dots + \kappa_d^{(l_1)} \ell_d^2 \right) \ell_{l_2} \\ &= \frac{1}{2} \omega^2 \left( \kappa_1^{(l_1)} \ell_1^2 \ell_{j_1} + \dots + \kappa_j^{(l_1)} \ell_j^3 + \dots + \kappa_d^{(l_1)} \ell_d^2 \ell_{l_2} \right).\end{aligned}$$

It follows that

$$(A.20) \quad \mathbb{E}[\omega^2 c_{l_1} \ell_{l_2}] = \frac{1}{2} \left( \kappa_1^{(l_1)} \mathbb{E}[\omega^2 \ell_1^2 \ell_{l_2}] + \dots + \kappa_{j_1}^{(l_1)} \mathbb{E}[\omega^2 \ell_{j_1}^3] + \dots + \kappa_d^{(l_1)} \mathbb{E}[\omega^2 \ell_d^2 \ell_{l_2}] \right).$$

Expanding the first type of expectation in Eq. (A.20) renders

$$\begin{aligned}\mathbb{E}[\omega^2 \ell_{j_2}^2 \ell_{l_2}] &= \mathbb{E} \left[ \left( 1 - \frac{\|\mathbf{x}\|^2}{s} \right) (r^2 u^{2/d} s_{j_2}^2) (r u^{1/d} s_{l_2}) \right] \\ &= \frac{dr^3}{d+3} \left( \mathbb{E}[s_{j_2}^2 s_{l_2}] - \frac{1}{s} \mathbb{E}[\|\mathbf{x}\|^2 s_{j_2}^2 s_{l_2}] \right) \\ &= \frac{-dr^3}{s(d+3)} \mathbb{E}[\|\mathbf{x}\|^2 s_{j_2}^2 s_{l_2}] \\ &= \frac{-dr^3}{s(d+3)} \left( \mathbb{E}[\|\boldsymbol{\ell}\|^2 s_{j_2}^2 s_{l_2}] + \mathbb{E}[\|\mathbf{c}\|^2 s_{j_2}^2 s_{l_2}] + \mathbb{E}[\|\mathbf{e}\|^2 s_{j_2}^2 s_{l_2}] + 2\mathbb{E}[\langle \boldsymbol{\ell} + \mathbf{c}, \mathbf{e} \rangle s_{j_2}^2 s_{l_2}] \right),\end{aligned}$$

for  $1 \leq j_2 \neq l_2 \leq d$ . By independence of the noise term, this reduces to

$$\mathbb{E}[\omega^2 \ell_{j_2}^2 \ell_{l_2}] = \frac{-dr^3}{s(d+3)} \left( \mathbb{E}[\|\boldsymbol{\ell}\|^2 s_{j_2}^2 s_{l_2}] + \mathbb{E}[\|\mathbf{c}\|^2 s_{j_2}^2 s_{l_2}] \right).$$

However, the remaining terms in this expression are also zero, because  $\mathbb{E}[s_{j_3}^2 s_{j_2}^2 s_{l_2}] = \mathbb{E}[s_{j_3}^4 s_{j_2}^2 s_{l_2}] = \mathbb{E}[s_{j_3}^2 s_{j_4}^2 s_{j_2}^2 s_{l_2}] = 0$ .

For the second type of term in Eq. (A.20), we have

$$\begin{aligned}\mathbb{E}[\omega^2 \ell_{l_1}^3] &= \frac{dr^3}{d+3} \left( \mathbb{E}[s_{l_1}^3] - \frac{1}{s} \mathbb{E}[\|\mathbf{x}\|^2 s_{l_1}^3] \right) \\ &= \frac{-dr^3}{s(d+3)} \left( \mathbb{E}[\|\boldsymbol{\ell}\|^2 s_{l_1}^3] + \mathbb{E}[\|\mathbf{c}\|^2 s_{l_1}^3] \right) \quad (\text{by independence of the noise terms}) \\ &= 0,\end{aligned}$$

## A.5. BOUNDING TERMS

---

because  $\mathbb{E}[s_{j_1}^2 s_{l_1}^3] = \mathbb{E}[s_{j_1}^4 s_{l_1}^3] = \mathbb{E}[s_{j_1}^2 s_{j_2}^2 s_{l_1}^3] = 0$ . Thus we conclude that  $\mathbb{E}[\omega^2 \mathbf{c}\ell^\top] = 0$ . Lastly, because the multiplication of  $\mathbf{c}\ell^\top$  by  $W_2^\top$  and  $W_1$  preserves its entries' zero expectation, it follows that  $\mathbb{E}[W_2^\top \omega^2 \mathbf{c}\ell^\top W_1] = 0$ .

Thus we have that

$$\begin{aligned} \|\hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c}\ell^\top W_1]\|_F &\leq \|\hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c}\ell^\top W_1 - \mathbb{E}[W_2^\top \omega^2 \mathbf{c}\ell^\top W_1]]\|_F + \|\mathbb{E}[W_2^\top \omega^2 \mathbf{c}\ell^\top W_1]\|_F \\ &= \|\hat{\mathbb{E}}[W_2^\top \omega^2 \mathbf{c}\ell^\top W_1 - \mathbb{E}[W_2^\top \omega^2 \mathbf{c}\ell^\top W_1]]\|_F \\ &\leq \frac{R_c R_\ell}{\sqrt{n}} (2 + \xi \sqrt{2}) \end{aligned}$$

with probability at least  $1 - e^{-\xi^2}$ . We label the corresponding probability constant  $\xi =: \xi_{cl}$ .

Plugging in  $R_c$  and  $R_\ell$ , we obtain

LEMMA A.5.6. *The Frobenius norm of the curvature-linear term satisfies*

$$\left\| W_2^\top \left( \frac{1}{n} CD_0 (LD_0)^\top \right) W_1 \right\|_F \leq \frac{K^{(+)} r^3}{2\sqrt{n}} (2 + \xi_{cl} \sqrt{2})$$

with probability at least  $1 - e^{-\xi_{cl}^2}$ . The quantity  $K^{(+)}$  is defined in Eq. (A.8).

**A.5.6. Noise Bounds.** To bound the Frobenius norm of the matrix  $W_i^\top (1/n) ED_0 (ED_0)^\top W_j$  where  $(i, j) \in \{(1, 1), (2, 2), (2, 1)\}$ , recall that the entries of  $Z := \alpha ED_0$  (for  $\alpha$  in Eq. (A.16)) are standard normal. Fortunately for our purposes, Kaslovsky and Meyer's original analysis for the analogous noise bound (in the standard local PCA case) is completely distribution-based, i.e., it depends only on the positions of nonzeros and the distribution of the matrix entries. Thus the results for the Singer and Wu case follow directly from making the necessary substitutions in Kaslovsky and Meyer's derived bounds, in particular, adjusting them to our definition of  $\alpha$ . Again, see their paper [57] for the derivation.

LEMMA A.5.7. *The following noise term satisfies*

$$\begin{aligned} \left\| W_1^\top \left( \frac{1}{n} ED_0 (ED_0)^\top \right) W_1 \right\|_F &\leq \sqrt{d} \left( \sigma^2 \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \right) \left( 1 - \frac{1}{n} \right)^{-1} \\ &\quad \times \left[ 1 + \frac{5}{2} \frac{1}{\sqrt{n}} (\sqrt{d} + \xi_{e_1} \sqrt{2}) \right] \end{aligned}$$

## A.5. BOUNDING TERMS

---

with probability greater than  $1 - 2e^{-\xi_{\mathbf{e}_1}^2}$ . The notation “ $\xi_{\mathbf{e}_1}$ ” indicates that  $W_1^\top ED_0(ED_0)^\top W_1$  only depends on the first  $d$  coordinates of the realizations of the random noise vector  $\mathbf{e} \in \mathbb{R}^m$ .

LEMMA A.5.8. *The following noise term satisfies*

$$\begin{aligned} \left\| W_2^\top \left( \frac{1}{n} ED_0(ED_0)^\top \right) W_2 \right\|_F &\leq \sqrt{m-d} \left( \sigma^2 \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \right) \left( 1 - \frac{1}{n} \right)^{-1} \\ &\quad \times \left[ 1 + \frac{5}{2} \frac{1}{\sqrt{n}} (\sqrt{m-d} + \xi_{\mathbf{e}_2} \sqrt{2}) \right] \end{aligned}$$

with probability greater than  $1 - 2e^{-\xi_{\mathbf{e}_2}^2}$ . The notation “ $\xi_{\mathbf{e}_2}$ ” indicates that  $W_2^\top ED_0(ED_0)^\top W_2$  only depends on the last  $m-d$  coordinates of the realizations of the random noise vector  $\mathbf{e} \in \mathbb{R}^m$ .

LEMMA A.5.9. *The following noise term satisfies*

$$\begin{aligned} \left\| W_2^\top \left( \frac{1}{n} ED_0(ED_0)^\top \right) W_1 \right\|_F &\leq \frac{\sqrt{d(m-d)}}{\sqrt{n}} \left( \sigma^2 \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \right) \left( 1 - \frac{1}{n} \right)^{-1} \\ &\quad \times \left( 1 + \frac{\sqrt{m-d} + \xi_{\mathbf{e}_2} \sqrt{2}}{\sqrt{n}} \right) \left[ 1 + \frac{6}{5} \frac{\xi_{\mathbf{e}_3}}{\sqrt{d(m-d)}} \right] \end{aligned}$$

with probability greater than  $1 - 2e^{-\xi_{\mathbf{e}_2}^2} - e^{-\xi_{\mathbf{e}_3}^2}$ . Here,  $\xi_{\mathbf{e}_3}$  is a probability constant that satisfies  $\xi_{\mathbf{e}_3} \leq 0.7\sqrt{d(m-d)}$ .

**A.5.7. Linear-Noise Bounds.** To bound  $\|W_i^\top (1/n) ED_0(LD_0)^\top W_1\|_F$  for  $i \in \{1, 2\}$ , define the matrices  $P_1 \in \mathbb{R}^{m \times d}$  and  $Q_1 \in \mathbb{R}^{d \times d}$  so that  $P_1 Q_1 = W_1$  with

$$P_1 = \begin{bmatrix} I_{d \times d} \\ 0_{(m-d) \times d} \end{bmatrix}.$$

Then it follows that  $Q_1$  is an orthogonal matrix, and we can write

$$\left\| W_1^\top \left( \frac{1}{n} ED_0(LD_0)^\top \right) U_1 \right\|_F = \frac{1}{n} \| (P_1 Q_1)^\top ED_0(LD_0)^\top (P_1 Q_1) \|_F = \frac{1}{n} \| P_1^\top ED_0(LD_0)^\top P_1 \|_F.$$

In order to simplify the term in the Frobenius norm, we expand  $P_1^\top LD_0$  in terms of its singular value decomposition:

$$\begin{aligned} P_1^\top LD_0 &= Q \Sigma V^\top \\ \Rightarrow \| P_1^\top \left( \frac{1}{n} ED_0(LD_0)^\top \right) P_1 \|_F &= \frac{1}{n} \| P_1^\top ED_0(Q \Sigma V^\top)^\top \|_F \end{aligned}$$

## A.5. BOUNDING TERMS

---

$$\begin{aligned}
&= \frac{1}{n} \|P_1^\top E D_0 V \Sigma Q^\top\|_F \\
&= \frac{1}{n} \|P_1^\top E D_0 V \Sigma\|_F.
\end{aligned}$$

Now,  $\Sigma$  contains the singular values of  $P_1^\top L D_0$ , which are equal to the square roots of the eigenvalues of  $P_1^\top L D_0 (P_1^\top L D_0)^\top$ . Since  $P_1$  is just a resizing matrix, these are equal to the square roots of the eigenvalues of  $L D_0 (L D_0)^\top$ . We can write

$$\begin{aligned}
\left\| P_1^\top \left( \frac{1}{n} E (L D_0)^\top \right) P_1 \right\|_F &= \frac{1}{n} \|P_1^\top E V \Sigma\|_F \\
&\leq \frac{1}{n} \|P_1^\top E V\|_F \sqrt{\lambda_{\max}(L D_0 (L D_0)^\top)} \\
&\leq \frac{1}{n} \|P_1^\top E V\|_F \sqrt{n \lambda_{\max}((1/n) L D_0 (L D_0)^\top)} \\
&\leq \frac{\sqrt{\lambda_1^W}}{\sqrt{n}} \|P_1^\top E V\|_F,
\end{aligned}$$

where  $\lambda_1^W := \lambda_{\max}((1/n) L D_0 (L D_0)^\top)$ , as before.

Despite the matrix  $V$ 's dependence on the realization of the matrix  $L D_0$ , the steps to bound  $\|P_1^\top E V\|_F$  follow identically to Kaslovsky and Meyer's analysis in the standard local PCA case. This is (again) a consequence of the standard normal distribution of the entries of  $\alpha E D_0$  (for  $\alpha$  in Eq. (A.16)), as well as the authors' approach of first bounding the (analogous) quantity for a fixed  $V$  followed by the removal of this conditioning.

We obtain the following bound:

**LEMMA A.5.10.** *The following noise-linear term satisfies*

$$\left\| W_1^\top \left( \frac{1}{n} E D_0 (L D_0)^\top \right) W_1 \right\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \frac{\sqrt{\lambda_{\text{bound}}^W(\xi_{\lambda_1^W})}}{\sqrt{n}} \left( d + \frac{6}{5} \xi_{e\ell} \right)$$

with probability greater than  $(1 - e^{-\xi_{e\ell}^2/2})(1 - de^{-\xi_{\lambda_1^W}^2})$ . Here,  $\xi_{e\ell}$  and  $\xi_{\lambda_1^W}$  are probability constants, and  $\lambda_{\text{bound}}^W$  is defined in Eq. (A.14).

The case that  $i = 2$  follows similarly:

## A.6. END RESULT

---

LEMMA A.5.11. *The following noise-linear term satisfies*

$$\left\| W_2^T \left( \frac{1}{n} ED_0 (LD_0)^T \right) W_1 \right\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \frac{\sqrt{\lambda_{\text{bound}}^W(\xi_{\lambda_1^W})}}{\sqrt{n}} \left( \sqrt{d(m-d)} + \frac{6}{5} \xi_{e\ell} \right)$$

*with probability greater than  $(1 - e^{-\xi_{e\ell}^2/2})(1 - de^{-\xi_{\lambda_1^W}^2})$ . Here,  $\xi_{e\ell}$  and  $\xi_{\lambda_1^W}$  are probability constants, and  $\lambda_{\text{bound}}^W$  is defined in Eq. (A.14).*

**A.5.8. Curvature-Noise Bounds.** Analogous to the original theorem, the bounds for  $\|W_2^T (1/n) CD_0 (ED_0)^T W_j\|_F$ ,  $j = 1, 2$ , can be obtained in a similar manner to the linear-noise bounds in Section A.5.7. We obtain the following lemmas:

LEMMA A.5.12. *The following curvature-noise term satisfies*

$$\left\| W_2^T \left( \frac{1}{n} CD_0 (ED_0)^T \right) W_1 \right\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \frac{\sqrt{\gamma_{\text{bound}}^W(\xi_c)}}{\sqrt{n}} \left( \sqrt{d(m-d)} + \frac{6}{5} \xi_{ce} \right)$$

*with probability greater than  $(1 - e^{-\xi_{ce}^2})(1 - 2e^{-\xi_c^2})$ . Here,  $\xi_{ce}$  and  $\xi_c$  are probability constants, and  $\gamma_{\text{bound}}^W$  is defined in Eq. (A.19).*

LEMMA A.5.13. *The following curvature-noise term satisfies*

$$\left\| W_2^T \left( \frac{1}{n} CD_0 (ED_0)^T \right) W_2 \right\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \frac{\sqrt{\gamma_{\text{bound}}^W(\xi_c)}}{\sqrt{n}} \left( m-d + \frac{6}{5} \xi_{ce} \right)$$

*with probability greater than  $(1 - e^{-\xi_{ce}^2})(1 - 2e^{-\xi_c^2})$ . Here,  $\xi_{ce}$  and  $\xi_c$  are probability constants, and  $\gamma_{\text{bound}}^W$  is defined in Eq. (A.19).*

## A.6. End Result

**A.6.1. Summary of Computed Bounds.** Let us consider the bounds computed in the previous section. We set  $\xi_{\lambda_d^W} = \xi_{\lambda_1^W} =: \xi_{\lambda^W}$  and

$$\xi_{c\ell} = \xi_{e\ell} = \xi_{ce} = \xi_{e_1} = \xi_{e_2} = \xi_{e_3} = \xi_c =: \xi.$$

We also multiply each bound by  $n$  (to remove the extra  $1/n$  factor), making the eigenvalue notation substitutions  $n\lambda_d^W =: \lambda_d$ ,  $n\lambda_{\text{bound}}^W =: \lambda_{\text{bound}}$ , and  $n\gamma_{\text{bound}}^W := \gamma_{\text{bound}}$  for convenience. This produces

## A.6. END RESULT

---

the equations

(A.21)

$$\begin{aligned} \lambda_{\text{bound}}(\xi) := & \frac{nr^2}{d+2} - \frac{n}{s} \left( \frac{r^2\sigma^2m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\ & + \xi \frac{5r^3\sqrt{n}}{2} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)} \end{aligned}$$

and

$$(A.22) \quad \gamma_{\text{bound}}(\xi) := n \left( \sum_{i,j=d+1}^m (\mathbb{E}[\omega^2 c_i c_j])^2 \right)^{1/2} + \frac{\sqrt{n} K^{(+)} (2 + \xi \sqrt{2})}{2},$$

where again,  $\mathbb{E}[\omega^2 c_i c_j]$  is defined in Eq. (A.18). It makes sense to then set the probability constant  $\xi_{\lambda W} =: \xi_\lambda$ .

This renders the following list of bounds:

- (1)  $\lambda_d > \frac{nr^2}{d+2} - \frac{n}{s} \left( \frac{r^2\sigma^2m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)$   
 $- \xi_\lambda r^3 \sqrt{2n} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)}$  with probability greater than  $1 - de^{-\xi_\lambda^2}$  and for  $\xi_\lambda$  satisfying Eq. (A.13) (by Lemma A.5.1),
- (2)  $\|W_2^\top C D_0 (C D_0)^\top W_2\|_F \leq n \left( \sum_{i,j=d+1}^m (\mathbb{E}[\omega^2 c_i c_j])^2 \right)^{1/2} + \frac{K^{(+)} \sqrt{n}(2 + \xi \sqrt{2})}{2} = \gamma_{\text{bound}}(\xi)$  with probability greater than  $1 - 2e^{-\xi^2}$ , where  $\mathbb{E}[\omega^2 c_i c_j]$  is given by Eq. (A.18) (by Lemma A.5.5),
- (3)  $\|W_2^\top C D_0 (L D_0)^\top W_1\|_F \leq \frac{K^{(+)} \sqrt{n} r^3}{2} (2 + \xi \sqrt{2})$  with probability at least  $1 - e^{-\xi^2}$  (by Lemma A.5.6),
- (4)  $\|W_1^\top E D_0 (E D_0)^\top W_1\|_F \leq \frac{n^2 \sigma^2 \sqrt{d}}{n-1} \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \left[ 1 + \frac{5}{2} \frac{1}{\sqrt{n}} (\sqrt{d} + \xi \sqrt{2}) \right]$  with probability greater than  $1 - 2e^{-\xi^2}$  (by Lemma A.5.7),
- (5)  $\|W_2^\top E D_0 (E D_0)^\top W_2\|_F \leq \frac{n^2 \sigma^2 \sqrt{m-d}}{n-1} \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \left[ 1 + \frac{5}{2} \frac{1}{\sqrt{n}} (\sqrt{m-d} + \xi \sqrt{2}) \right]$  with probability greater than  $1 - 2e^{-\xi^2}$  (by Lemma A.5.8),
- (6)  $\|W_2^\top E D_0 (E D_0)^\top W_1\|_F \leq \frac{n^{3/2} \sigma^2 \sqrt{d(m-d)}}{n-1} \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \left( 1 + \frac{\sqrt{m-d} + \xi \sqrt{2}}{\sqrt{n}} \right) \left[ 1 + \frac{6}{5} \frac{\xi}{\sqrt{d(m-d)}} \right]$  with probability greater than  $1 - 2e^{-\xi^2} - e^{-\xi^2}$  and for  $\xi \leq 0.7 \sqrt{d(m-d)}$  (by Lemma A.5.9),

## A.6. END RESULT

---

- (7)  $\|W_1^T ED_0(LD_0)^T W_1\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\lambda_{\text{bound}}(\xi_\lambda)} \left( d + \frac{6}{5} \xi \right)$  with probability greater than  $(1 - e^{-\xi^2/2})(1 - de^{-\xi_\lambda^2})$  (by Lemma A.5.10),
- (8)  $\|W_2^T ED_0(LD_0)^T W_1\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\lambda_{\text{bound}}(\xi_\lambda)} \left( \sqrt{d(m-d)} + \frac{6}{5} \xi \right)$  with probability greater than  $(1 - e^{-\xi^2/2})(1 - de^{-\xi_\lambda^2})$  (by Lemma A.5.11),
- (9)  $\|W_2^T CD_0(ED_0)^T W_1\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\gamma_{\text{bound}}(\xi)} \left( \sqrt{d(m-d)} + \frac{6}{5} \xi \right)$  with probability greater than  $(1 - e^{-\xi^2})(1 - 2e^{-\xi^2})$  (by Lemma A.5.12),
- (10)  $\|W_2^T CD_0(ED_0)^T W_2\|_F \leq \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\gamma_{\text{bound}}(\xi)} \left( m - d + \frac{6}{5} \xi \right)$  (by Lemma A.5.13).

**A.6.2. Putting It All Together.** We restate Theorem A.1.1 with our notation for easy reference:

THEOREM A.6.1 (Restated theorem: Davis and Kahan [25], Stewart [89] ). Set

$$\delta := \lambda_d - \|W_1^T(n\Theta)W_1\|_F - \|W_2^T(n\Theta)W_2\|_F.$$

Define the following conditions:

- Condition 1:  $\delta > 0$ ,
- Condition 2:  $\|W_1^T(n\Theta)W_2\|_F \|W_2^T(n\Theta)W_1\|_F < \frac{1}{4}\delta^2$ .

If both conditions hold, then

$$\|WW^T - \widehat{W}\widehat{W}^T\|_F \leq 2\sqrt{2} \frac{\|W_2^T(n\Theta)W_1\|_F}{\delta}.$$

From the triangle inequality and the equalities  $W_1^T C = 0$  and  $W_2^T L = 0$ , we have

$$\begin{aligned} \|W_1^T(n\Theta)W_1\|_F &\leq 2\|W_1^T LD_0(ED_0)W_1\|_F + \|W_1^T ED_0(ED_0)W_1\|_F, \\ \|W_2^T(n\Theta)W_2\|_F &\leq \|W_2^T CD_0(CD_0)W_2\|_F + 2\|W_2^T CD_0(ED_0)W_2\|_F + \|W_2^T ED_0(ED_0)W_2\|_F, \\ \|W_2^T(n\Theta)W_1\|_F &\leq \|W_2^T CD_0(LD_0)W_1\|_F + \|W_2^T ED_0(LD_0)W_1\|_F \\ &\quad + \|W_2^T CD_0(ED_0)W_1\|_F + \|W_2^T ED_0(ED_0)W_1\|_F. \end{aligned}$$

It now is simply a matter of plugging the computed bounds into Theorem A.6.1.

## A.6. END RESULT

---

THEOREM A.6.2 (Tangent bound, modified for Singer and Wu's method of local PCA). *Let  $s$  be the scaling factor defined in Eq. (A.9), and let  $\xi_\lambda$  satisfying Eq. (A.13) and  $\xi \leq 0.7\sqrt{d(m-d)}$  be probability constants with  $n > 4(\max(\sqrt{d}, \sqrt{m-d}) + \xi)$ . Assume the setup defined above, in particular: (i) the eigendecompositions in Eq. (A.5) and Eq. (A.6); (ii) the curvature constants defined in Eq. (A.8), Eq. (A.12), Eq. (A.15), and Eq. (A.17); and (iii) the eigenvalue bounds in Eq. (A.21) and Eq. (A.22) (the latter of which uses Eq. (A.18)). Further define*

$$\begin{aligned} \delta := & \frac{nr^2}{d+2} - \frac{n}{s} \left( \frac{r^2\sigma^2 m}{d+2} + \frac{r^4}{d+4} + \frac{r^6}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right) \\ & - \xi_\lambda r^3 \sqrt{2n} \sqrt{\frac{1}{d+2} - \frac{1}{s} \left( \frac{\sigma^2 m}{d+2} + \frac{r^2}{d+4} + \frac{r^4}{4(d+6)} \left[ \frac{3(d+4)}{(d+2)(d+4)} (K^{(+)})^2 + \frac{d-1}{d+2} \sum_{i=d+1}^m K_{j_1 j_2}^{ii} \right] \right)} \\ & - 2\sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\lambda_{\text{bound}}(\xi_\lambda)} \left( d + \frac{6}{5}\xi \right) \\ & - \frac{n^2 \sigma^2 \sqrt{d}}{n-1} \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \left[ 1 + \frac{5}{2} \frac{1}{\sqrt{n}} (\sqrt{d} + \xi \sqrt{2}) \right] \\ & - \gamma_{\text{bound}}(\xi) - 2\sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\gamma_{\text{bound}}(\xi)} \left( m - d + \frac{6}{5}\xi \right) \\ & - \frac{n^2 \sigma^2 \sqrt{m-d}}{n-1} \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \left[ 1 + \frac{5}{2} \frac{1}{\sqrt{n}} (\sqrt{m-d} + \xi \sqrt{2}) \right], \end{aligned}$$

and

$$\begin{aligned} \beta := & \frac{K^{(+)} \sqrt{n} r^3}{2} (2 + \xi \sqrt{2}) + \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\lambda_{\text{bound}}(\xi_\lambda)} \left( \sqrt{d(m-d)} + \frac{6}{5}\xi \right) \\ & + \sigma \sqrt{1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right)} \sqrt{\gamma_{\text{bound}}(\xi)} \left( \sqrt{d(m-d)} + \frac{6}{5}\xi \right) \\ & + \frac{n^{3/2} \sigma^2 \sqrt{d(m-d)}}{n-1} \left[ 1 - \frac{1}{s} \left( \frac{r^2 d}{d+2} + \gamma r^4 + \sigma^2 m \right) \right] \left( 1 + \frac{\sqrt{m-d} + \xi \sqrt{2}}{\sqrt{n}} \right) \left[ 1 + \frac{6}{5} \frac{\xi}{\sqrt{d(m-d)}} \right]. \end{aligned}$$

If the two conditions

- $\delta > 0$ , and
- $\beta < \frac{1}{2}\delta$

## A.7. ADJUSTING THE SCALING FACTOR

---

hold, then the angle between the true and estimated tangent plane, computed using Singer and Wu's local PCA algorithm, satisfies

$$(A.23) \quad \|WW^\top - \widehat{W}\widehat{W}^\top\|_F \leq \frac{2\sqrt{2}\beta}{\delta}$$

with probability greater than  $1 - 2de^{-\xi_\lambda^2} - 9e^{-\xi^2}$ .

### A.7. Adjusting the Scaling Factor

As a final consideration, we investigate how the bound in Theorem A.6.2 is affected when a smaller scaling factor  $s$  is used, i.e., when we use a tighter upper bound on  $\|\mathbf{x}\|^2$  than the one in Eq. (A.9) (recall that this is the case in the LPCA-SRC algorithm). By the quantities in Theorem A.6.2, making  $s$  smaller will cause  $\beta$  and the subtracted terms in the expression for  $\delta$  to decrease. However, the first term in the expression for  $\delta$  (the lower bound on  $\lambda_d$ ) also decreases when  $s$  is made smaller, and so the overall behavior of the tangent bound is not immediately obvious.

We expect that, in general, the tangent bound will decrease as  $s$  decreases to  $\max_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_2^2 + \epsilon$ , for a fixed set of neighbor samples  $\mathcal{X}$  and a small but non-negligible constant  $\epsilon > 0$ . In the case that  $d = 1$  (as in our tangent vector distance theorem in Chapter 6),  $\delta$  can be written informally as

$$\begin{aligned} \delta \approx n \left[ r^2 \left[ \frac{1}{3} - \frac{1}{s} \left( \frac{\sigma^2 m}{3} + \frac{1}{5} + \frac{3r^2 \|\boldsymbol{\kappa}\|_2^2}{84} \right) \right] - r^4 \|\boldsymbol{\kappa}\|_2 \left( \frac{1}{20} - \frac{1}{s} \left( \frac{\sigma^2}{20} + \frac{r^2}{7} + \frac{r^2 \|\boldsymbol{\kappa}\|_2^2}{144} \right) \right) \right. \\ \left. - \sigma^2 \left[ 1 - \frac{1}{s} \left( \frac{r^2}{3} + \frac{r^4 \|\boldsymbol{\kappa}\|_2^2}{10} + \sigma^2 m \right) \right] (1 + \sqrt{m-1}) \right], \end{aligned}$$

where  $\boldsymbol{\kappa} := [\kappa_1^{(2)}, \dots, \kappa_1^{(m)}]^\top \in \mathbb{R}^{m-1}$ . Thus we can (heuristically) confirm that the bound in Eq. (A.23) decreases as  $s$  decreases by checking to see if

$$\begin{aligned} \frac{1}{n} \delta'(s) \approx r^2 \left[ \frac{1}{3} + \frac{1}{s^2} \left( \frac{\sigma^2 m}{3} + \frac{1}{5} + \frac{3r^2 \|\boldsymbol{\kappa}\|_2^2}{84} \right) \right] - r^4 \|\boldsymbol{\kappa}\|_2 \left( \frac{1}{20} + \frac{1}{s^2} \left( \frac{\sigma^2}{20} + \frac{r^2}{7} + \frac{r^2 \|\boldsymbol{\kappa}\|_2^2}{144} \right) \right) \\ - \sigma^2 \left[ 1 + \frac{1}{s^2} \left( \frac{r^2}{3} + \frac{r^4 \|\boldsymbol{\kappa}\|_2^2}{10} + \sigma^2 m \right) \right] (1 + \sqrt{m-1}) \end{aligned}$$

## A.7. ADJUSTING THE SCALING FACTOR

---

is non-negative. However, we note that this is not sufficient to reject this hypothesis, and that doing so would require showing that

$$g'(s) := \frac{\partial}{\partial s} \left( \frac{\beta(s)}{\delta(s)} \right)$$

is less than 0. Unfortunately, the expansion of  $g'(s)$ , even in the case that  $d = 1$ , is quite complicated.

## Bibliography

- [1] ASIF, M., AND ROMBERG, J.  $\ell_1$  homotopy: A MATLAB toolbox for homotopy algorithms in  $\ell_1$ -norm minimization problems. <http://users.ece.gatech.edu/~sasif/homotopy/>, accessed 31.3.2015. 2009–2013.
- [2] AT&T LABORATORIES CAMBRIDGE. The database of faces. <http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>, accessed 26.3.2016. 1992–1994.
- [3] BELHUMEUR, P. N., AND KRIEGMAN, D. J. What is the set of images of an object under all possible lighting conditions? In *1996 IEEE Conference on Computer Vision and Pattern Recognition* (June 1996), pp. 270–277.
- [4] BELKIN, M., AND NIYOGI, P. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computing* 15, 6 (2003), 1373–1396.
- [5] BENTLEY, J. L. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18, 9 (1975), 509–517.
- [6] BEYREUTHER, M., CARNIEL, R., AND WASSERMANN, J. Automatic earthquake detection and classification with continuous hidden Markov models: a possible tool for monitoring Las Candas caldera in Tenerife. *IOP Conference Series: Earth and Environmental Science* 3, 1 (2008), 12–28.
- [7] BOSER, B. E., GUYON, I. M., AND VAPNIK, V. N. A training algorithm for optimal margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory* (1992), COLT '92, ACM, pp. 144–152.
- [8] BRUCKSTEIN, A. M., DONOHO, D. L., AND ELAD, M. From sparse solutions of systems of equations to sparse modeling of signals and images. *SIAM Rev.* 51, 1 (2009), 34–81.
- [9] CAI, T. T., WANG, L., AND XU, G. New bounds for restricted isometry constants. *IEEE Trans. Inform. Theory* 56, 9 (2010), 4388–4394.
- [10] CANDÈS, E. J., AND PLAN, Y. Near-ideal model selection by  $\ell_1$  minimization. *Ann. Statist.* 37, 5A (2009), 2145–2177.
- [11] CANDÈS, E. J., AND PLAN, Y. A probabilistic and RIPless theory of compressed sensing. *IEEE Trans. Inform. Theory* 57, 11 (2011), 7235–7254.
- [12] CANDÈS, E. J., ROMBERG, J., AND TAO, T. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inform. Theory* 52, 2 (2006), 489–509.
- [13] CANDÈS, E. J., ROMBERG, J. K., AND TAO, T. Stable signal recovery from incomplete and inaccurate measurements. *Comm. Pure Appl. Math.* 59, 8 (2006), 1207–1223.

- 
- [14] CANDÈS, E. J., AND TAO, T. Decoding by linear programming. *IEEE Trans. Inform. Theory* 51, 12 (2005), 4203–4215.
- [15] CANDÈS, E. J., AND TAO, T. Near-optimal signal recovery from random projections: universal encoding strategies? *IEEE Trans. Inform. Theory* 52, 12 (2006), 5406–5425.
- [16] CERUTI, C., BASSIS, S., ROZZA, A., LOMBARDI, G., CASIRAGHI, E., AND CAMPADELLI, P. DANCo: An intrinsic dimensionality estimator exploiting angle and norm concentration. *Pattern Recognition* 47, 8 (2014), 2569 – 2581.
- [17] CEVIKALP, H., YAVUZ, H. S., CAY, M. A., AND BARKANA, A. Two-dimensional subspace classifiers for face recognition. *Neurocomputing* 72, 46 (2009), 1111 – 1120.
- [18] CHANG, J.-M., AND KIRBY, M. Face recognition under varying viewing conditions with subspace distance. In *International Conference on Artificial Intelligence and Pattern Recognition (AIPR-09)* (2009), pp. 16–23.
- [19] CHANG, J.-M., AND PACHECO, J. I. Classifying handwritten digits on the Grassmann manifold. In *Proceedings of The International Conference on Image Processing, Computer Vision, & Pattern Recognition* (2011), vol. 1, pp. 36–41.
- [20] CHANG, Y., HU, C., AND TURK, M. Manifold of facial expression. In *Proceedings of the IEEE International Workshop on Analysis and Modeling of Faces and Gestures* (2003), vol. 24, pp. 605–614.
- [21] CHENG, B., YANG, J., YAN, S., FU, Y., AND HUANG, T. S. Learning with  $l_1$ -graph for image analysis. *IEEE Trans. Image Process.* 19, 4 (2010), 858–866.
- [22] CORTES, C., AND VAPNIK, V. Support-vector networks. *Machine Learning* 20, 3 (1995), 273–297.
- [23] COVER, T., AND HART, P. Nearest neighbor pattern classification. *IEEE Trans. Inform. Theory* 13, 1 (1967), 21–27.
- [24] DALAL, N., AND TRIGGS, B. Histograms of oriented gradients for human detection. In *2005 IEEE Conference on Computer Vision and Pattern Recognition* (June 2005), vol. 1, pp. 886–893.
- [25] DAVIS, C., AND KAHAN, W. M. The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.* 7 (1970), 1–46.
- [26] DENG, W., HU, J., AND GUO, J. Extended SRC: Undersampled face recognition via intra-class variant dictionary. *IEEE Trans. Pattern Anal. Mach. Intell.* 34, 9 (2012), 1864–1870.
- [27] DONOHO, D. L. Neighborly polytopes and sparse solutions of underdetermined linear equations. Tech. rep., 2005.
- [28] DONOHO, D. L. Compressed sensing. *IEEE Trans. Inform. Theory* 52, 4 (2006), 1289–1306.
- [29] DONOHO, D. L. For most large underdetermined systems of linear equations the minimal  $l_1$ -norm solution is also the sparsest solution. *Comm. Pure Appl. Math.* 59, 6 (2006), 797–829.
- [30] DONOHO, D. L., AND ELAD, M. Optimally sparse representation in general (nonorthogonal) dictionaries via  $\ell^1$  minimization. *Proc. Natl. Acad. Sci. USA* 100, 5 (2003), 2197–2202.

- 
- [31] DONOHO, D. L., ELAD, M., AND TEMLYAKOV, V. N. Stable recovery of sparse overcomplete representations in the presence of noise. *IEEE Trans. Inform. Theory* 52, 1 (2006), 6–18.
- [32] DONOHO, D. L., AND TANNER, J. Sparse nonnegative solution of underdetermined linear equations by linear programming. *Proc. Natl. Acad. Sci. USA* 102, 27 (2005), 9446–9451.
- [33] DONOHO, D. L., AND TSAIG, Y. Fast solution of  $l_1$ -norm minimization problems when the solution may be sparse. *IEEE Trans. Inform. Theory* 54, 11 (2008), 4789–4812.
- [34] DUARTE, M. F., DAVENPORT, M. A., TAKBAR, D., LASKA, J. N., SUN, T., KELLY, K. F., AND BARANIUK, R. G. Single-pixel imaging via compressive sampling. *IEEE Signal Processing Magazine* 25, 2 (March 2008), 83–91.
- [35] EDELMAN, A. Eigenvalues and condition numbers of random matrices. *SIAM J. Matrix Anal. Appl.* 9, 4 (1988), 543–560.
- [36] ELDAR, Y. C., KUPPINGER, P., AND BÖLCSKEI, H. Block-sparse signals: uncertainty relations and efficient recovery. *IEEE Trans. Signal Process.* 58, 6 (2010), 3042–3054.
- [37] ELDAR, Y. C., AND MISHALI, M. Robust recovery of signals from a structured union of subspaces. *IEEE Trans. Inform. Theory* 55, 11 (2009), 5302–5316.
- [38] ELDÉN, L. *Matrix methods in data mining and pattern recognition*, vol. 4 of *Fundamentals of Algorithms*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2007.
- [39] ELHAMIFAR, E., AND VIDAL, R. Sparse subspace clustering. In *2009 IEEE Conference on Computer Vision and Pattern Recognition* (June 2009), pp. 2790–2797.
- [40] ELHAMIFAR, E., AND VIDAL, R. Robust classification using structured sparse representation. In *2011 IEEE Conference on Computer Vision and Pattern Recognition* (June 2011), pp. 1873–1879.
- [41] ELHAMIFAR, E., AND VIDAL, R. Block-sparse recovery via convex optimization. *IEEE Trans. Signal Process.* 60, 8 (2012), 4094–4107.
- [42] FISHER, R. A. The use of multiple measurements in taxonomic problems. *Annals of Eugenics* 7, 2 (1936), 179–188.
- [43] FRIEDMAN, J., HASTIE, T., AND TIBSHIRANI, R. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33, 1 (2010), 1–22.
- [44] GEORGHIADES, A. S., BELHUMEUR, P. N., AND KRIEGMAN, D. J. From few to many: illumination cone models for face recognition under variable lighting and pose. *IEEE Trans. Pattern Anal. Mach. Intell.* 23, 6 (2001), 643–660.
- [45] GIAQUINTA, M., AND MODICA, G. *Mathematical Analysis: An Introduction to Functions of Several Variables*. Birkhäuser Boston, Inc., Boston, MA, 2009. Translated and revised from the 2005 Italian original.
- [46] GRIBONVAL, R., AND NIELSEN, M. Sparse representations in unions of bases. *IEEE Trans. Inform. Theory* 49, 12 (2003), 3320–3325.

- 
- [47] HAND, D. J., AND HENLEY, W. E. Statistical classification methods in consumer credit scoring: a review. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 160, 3 (1997), 523–541.
- [48] HASTIE, T., TIBSHIRANI, R., AND FRIEDMAN, J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, second ed. Springer Series in Statistics. Springer, New York, 2009.
- [49] HASTIE, T., TIBSHIRANI, R., AND WAINWRIGHT, M. *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press, Taylor & Francis, 2015.
- [50] HE, X., YAN, S., HU, Y., NIYOGI, P., AND ZHANG, H.-J. Face recognition using Laplacianfaces. *IEEE Trans. Pattern Anal. Mach. Intell.* 27, 3 (2005), 328–340.
- [51] HERRMANN, F. J., FRIEDLANDER, M. P., AND YILMAZ, O. Fighting the curse of dimensionality: Compressive sensing in exploration seismology. *IEEE Signal Processing Magazine* 29, 3 (May 2012), 88–100.
- [52] HO, J., XIE, Y., AND VEMURI, B. C. On a nonlinear generalization of sparse coding and dictionary learning. In *ICML (3)* (2013), vol. 28 of *JMLR Proceedings*, JMLR.org, pp. 1480–1488.
- [53] HUI, K.-H., LI, C.-L., AND ZHANG, L. Sparse neighbor representation for classification. *Pattern Recognition Letters* 33, 5 (2012), 661 – 669.
- [54] HULL, J. J. A database for handwritten text recognition research. *IEEE Trans. Pattern Anal. Mach. Intell.* 16, 5 (1994), 550–554.
- [55] KANG, C., LIAO, S., XIANG, S., AND PAN, C. Kernel sparse representation with pixel-level and region-level local feature kernels for face recognition. *Neurocomputing* 133 (2014), 141 – 152.
- [56] KASLOVSKY, D. Matlab code for “Non-Asymptotic Analysis of Tangent Space Perturbation”. <http://danielkaslovsky.com/code>, accessed 18.6.15. 2012.
- [57] KASLOVSKY, D. N., AND MEYER, F. G. Non-asymptotic analysis of tangent space perturbation. *Inf. Inference* 3, 2 (2014), 134–187.
- [58] LECUN, Y., BOTTOU, L., BENGIO, Y., AND HAFFNER, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE* 86, 11 (Nov 1998), 2278–2324.
- [59] LEE, K.-C., HO, J., AND KRIEGMAN, D. Acquiring linear subspaces for face recognition under variable lighting. *IEEE Trans. Pattern Anal. Mach. Intell.* 27, 5 (2005), 684–698.
- [60] LI, C.-G., GUO, J., AND ZHANG, H.-G. Local sparse representation based classification. In *2010 20th International Conference on Pattern Recognition* (Aug 2010), pp. 649–652.
- [61] LI, Q., LI, T., ZHU, S., AND KAMBHAMETTU, C. Improving medical/biological data classification performance by wavelet preprocessing. In *2002 IEEE International Conference on Data Mining* (2002), pp. 657–660.
- [62] LITTLE, A. V., MAGGIONI, M., AND ROSASCO, L. Multiscale geometric methods for data sets I: Multiscale SVD, noise and curvature. *Appl. Comput. Harmon. Anal.* (2016), in press.
- [63] LOWE, D. G. Object recognition from local scale-invariant features. In *Proceedings of the International Conference on Computer Vision* (1999), vol. 2, IEEE Computer Society, pp. 1150–1157.

- 
- [64] LUSTIG, M., DONOHO, D. L., SANTOS, J. M., AND PAULY, J. M. Compressed sensing MRI. *IEEE Signal Processing Magazine* 25, 2 (March 2008), 72–82.
- [65] MARTINEZ, A., AND BENAVENTE, R. The AR face database. Tech. Rep. 24, Computer Vision Center, June 1998.
- [66] MENG, H., PEARS, N., AND BAILEY, C. Human action classification using SVM\_2K classifier on motion features. In *Multimedia Content Representation, Classification and Security*, B. Gunsel, A. Jain, A. Tekalp, and B. Sankur, Eds., vol. 4105 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2006, pp. 458–465.
- [67] MERKWIRTH, C., PARLITZ, U., AND LAUTERBORN, W. Fast nearest-neighbor searching for nonlinear signal processing. *Phys. Rev. E* 62 (2000), 2089–2097.
- [68] MERKWIRTH, C., PARLITZ, U., WEDEKIND, I., ENGSTER, D., AND LAUTERBORN, W. TSTOOL homepage. <http://www.physik3.gwdg.de/tstool/index.html>, accessed 6.2.15. 2009.
- [69] NIYOGI, P., SMALE, S., AND WEINBERGER, S. Finding the homology of submanifolds with high confidence from random samples. *Discrete Comput. Geom.* 39, 1-3 (2008), 419–441.
- [70] OJALA, T., PIETIKAINEN, M., AND HARWOOD, D. Performance evaluation of texture measures with classification based on Kullback discrimination of distributions. In *Pattern Recognition, 1994. Conference A: Proceedings of the 12th IAPR International Conference on Computer Vision and Image Processing* (1994), vol. 1, pp. 582–585.
- [71] PATEL, R., RATHOD, N., AND SHAH, A. Comparative analysis of face recognition approaches: A survey. *International Journal of Computer Applications* 57, 17 (Nov 2012), 50–69.
- [72] PENNEBAKER, W. B., AND MITCHELL, J. L. *JPEG: Still Image Data Compression Standard*, 1st ed. Kluwer Academic Publishers, Norwell, MA, USA, 1992.
- [73] QIAO, L., CHEN, S., AND TAN, X. Sparsity preserving projections with applications to face recognition. *Pattern Recognition* 43, 1 (2010), 331–341.
- [74] RADHAKRISHNA RAO, C. The utilization of multiple measurements in problems of biological classification. *J. Roy. Statist. Soc. Ser. B.* 10 (1948), 159–193.
- [75] RAGEL, A., AND CRÉMILLEUX, B. MVC—a preprocessing method to deal with missing values. *Knowledge-Based Systems* 12, 5-6 (1999), 285 – 291.
- [76] RAY, A., SHUKLA, K., AGGARWAL, L., SHARMA, N., PRADHAN, S., AND SHARMA, S. Segmentation and classification of medical images using texture-primitive features: Application of BAM-type artificial neural network. *Journal of Medical Physics* 33, 3 (2008), 119–126.
- [77] ROSENFIELD, M. In praise of the Gram matrix. In *The mathematics of Paul Erdős, II*, vol. 14 of *Algorithms Combin.* Springer, Berlin, 1997, pp. 318–323.
- [78] ROWEIS, S. T., AND SAUL, L. K. Nonlinear dimensionality reduction by locally linear embedding. *Science* 290 (2000), 2323–2326.

- 
- [79] SAITO, N. *Local Feature Extraction and Its Applications Using a Library of Bases*. PhD thesis, Department of Mathematics, Yale University, Dec 1994.
- [80] SAITO, N., AND COIFMAN, R. R. Local discriminant bases and their applications. *J. Math. Imaging Vis.* 5, 4 (1995), 337–358. Invited paper.
- [81] SEUNG, H. S., AND LEE, D. D. The manifold ways of perception. *Science* 290, 5500 (2000), 2268–2269.
- [82] SHANNON, C. E. Communication in the presence of noise. *Proc. I.R.E.* 37 (1949), 10–21.
- [83] SHARON, Y., WRIGHT, J., AND MA, Y. Minimum sum of distances estimator: Robustness and stability. In *2009 American Control Conference* (June 2009), pp. 524–530.
- [84] SHAWE-TAYLOR, J., AND CRISTIANINI, N. Estimating the moments of a random vector with applications. In *Proceedings of GRETSI 2003 Conference*, vol. I. 2003, pp. 47–52. Invited talk.
- [85] SIMARD, P. Y., LECUN, Y. A., DENKER, J. S., AND VICTORRI, B. *Neural Networks: Tricks of the Trade: Second Edition*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012, ch. Transformation Invariance in Pattern Recognition – Tangent Distance and Tangent Propagation, pp. 235–269.
- [86] SIMON, N., FRIEDMAN, J., HASTIE, T., AND TIBSHIRANI, R. Regularization paths for Cox’s proportional hazards model via coordinate descent. *Journal of Statistical Software* 39, 5 (2011), 1–13.
- [87] SINGER, A., AND WU, H.-T. Vector diffusion maps and the connection Laplacian. *Comm. Pure Appl. Math.* 65, 8 (2012), 1067–1144.
- [88] SPRECHMANN, P., RAMÍREZ, I., SAPIRO, G., AND ELDAR, Y. C. C-HiLasso: a collaborative hierarchical sparse modeling framework. *IEEE Trans. Signal Process.* 59, 9 (2011), 4183–4198.
- [89] STEWART, G. W., AND SUN, J. G. *Matrix Perturbation Theory*. Computer Science and Scientific Computing. Academic Press, Inc., Boston, MA, 1990.
- [90] SUO, Y., DAO, M., TRAN, T., MOUSAVI, H., SRINIVAS, U., AND MONGA, V. Group structured dirty dictionary learning for classification. In *2014 IEEE International Conference on Image Processing (ICIP)* (Oct 2014), pp. 150–154.
- [91] TIMOFTE, R., AND GOOL, L. V. Adaptive and weighted collaborative representations for image classification. *Pattern Recognition Letters* 43 (2014), 127 – 135.
- [92] TROPP, J. A. On the conditioning of random subdictionaries. *Appl. Comput. Harmon. Anal.* 25, 1 (2008), 1–24.
- [93] TROPP, J. A. User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.* 12, 4 (2012), 389–434.
- [94] TYAGI, H., VURAL, E., AND FROSSARD, P. Tangent space estimation for smooth embeddings of Riemannian manifolds. *Inf. Inference* 2, 1 (2013), 69–114.
- [95] VAN DEN BERG, E., AND FRIEDLANDER, M. P. SPGL1: A solver for large-scale sparse reconstruction, June 2007. Version 1.9, April 2015 (accessed 12.4.2016).

- 
- [96] VAN DEN BERG, E., AND FRIEDLANDER, M. P. Probing the Pareto frontier for basis pursuit solutions. *SIAM Journal on Scientific Computing* 31, 2 (2008), 890–912.
- [97] VERSHYNIN, R. Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing*. Cambridge Univ. Press, Cambridge, 2012, pp. 210–268.
- [98] WANG, Z., YANG, J., NASRABADI, N., AND HUANG, T. A max-margin perspective on sparse representation-based classification. In *2013 IEEE International Conference on Computer Vision* (Dec 2013), pp. 1217–1224.
- [99] WAQAS, J., YI, Z., AND ZHANG, L. Collaborative neighbor representation based classification using  $l_2$ -minimization approach. *Pattern Recognition Letters* 34, 2 (2013), 201 – 208.
- [100] WEAVER, C., AND SAITO, N. Improving sparse representation-based classification using local principal component analysis. *CoRR abs/1204.2358* (2016). Submitted for publication.
- [101] WEI, C.-P., CHAO, Y.-W., YEH, Y.-R., AND WANG, Y.-C. F. Locality-sensitive dictionary learning for sparse representation based classification. *Pattern Recognition* 46, 5 (2013), 1277–1287.
- [102] WRIGHT, J., AND MA, Y. Dense error correction via  $\ell^1$ -minimization. *IEEE Trans. Inform. Theory* 56, 7 (2010), 3540–3560.
- [103] WRIGHT, J., MA, Y., MAIRAL, J., SAPIRO, G., HUANG, T. S., AND YAN, S. Sparse representation for computer vision and pattern recognition. *Proceedings of the IEEE* 98, 6 (June 2010), 1031–1044.
- [104] WRIGHT, J., YANG, A. Y., GANESH, A., SASTRY, S. S., AND MA, Y. Robust face recognition via sparse representation. *IEEE Trans. Pattern Anal. Mach. Intell.* 31, 2 (2009), 210–227.
- [105] XIAOYAN, Z., HOUJUN, W., AND ZHIJIAN, D. Wireless sensor networks based on compressed sensing. In *3rd IEEE International Conference on Computer Science and Information Technology (ICCSIT)* (July 2010), vol. 9, pp. 90–92.
- [106] XU, Y., LI, X., YANG, J., AND ZHANG, D. Integrate the original face image and its mirror image for face recognition. *Neurocomputing* 131 (2014), 191–199.
- [107] YANG, A. Y., SASTRY, S. S., GANESH, A., AND MA, Y. Fast  $\ell^1$ -minimization algorithms and an application in robust face recognition: A review. In *2010 17th IEEE International Conference on Image Processing* (Sept 2010), pp. 1849–1852.
- [108] YANG, J., WANG, J., AND HUANG, T. Learning the sparse representation for classification. In *2011 IEEE International Conference on Multimedia and Expo (ICME)* (July 2011), pp. 1–6.
- [109] YANG, J., ZHANG, L., XU, Y., AND YANG, J.-Y. Beyond sparsity: The role of L1-optimizer in pattern classification. *Pattern Recognition* 45, 3 (2012), 1104–1118.
- [110] YANG, J., ZHU, K., AND ZHONG, N. Local tangent distances for classification problems. In *2012 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology (WI-IAT)* (Dec 2012), vol. 1, pp. 396–401.

- 
- [111] YIN, J., LIU, Z., JIN, Z., AND YANG, W. Kernel sparse representation based classification. *Neurocomputing* 77, 1 (2012), 120 – 128.
  - [112] YUAN, M., AND LIN, Y. Model selection and estimation in regression with grouped variables. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 68, 1 (2006), 49–67.
  - [113] ZHANG, H., WANG, F., CHEN, Y., ZHANG, W., WANG, K., AND LIU, J. Sample pair based sparse representation classification for face recognition. *Expert Systems with Applications* 45 (2016), 352 – 358.
  - [114] ZHANG, L., YANG, M., AND FENG, X. Sparse representation or collaborative representation: Which helps face recognition? In *Proceedings of the 2011 International Conference on Computer Vision* (Dec 2011), IEEE Computer Society, pp. 471–478.
  - [115] ZHANG, L., YANG, M., FENG, X., MA, Y., AND ZHANG, D. Collaborative representation based classification for face recognition. *CoRR abs/1204.2358* (2012).
  - [116] ZHANG, L., ZHOU, W.-D., CHANG, P.-C., LIU, J., YAN, Z., WANG, T., AND LI, F.-Z. Kernel sparse representation-based classifier. *IEEE Trans. Signal Process.* 60, 4 (2012), 1684–1695.
  - [117] ZHANG, X., PHAM, D.-S., VENKATESH, S., LIU, W., AND PHUNG, D. Mixed-norm sparse representation for multiview face recognition. *Pattern Recognition* 48, 9 (2015), 2935 – 2946.
  - [118] ZHOU, Y., GAO, J., AND BARNER, K. E. An enhanced sparse representation strategy for signal classification. In *Compressive Sensing*, vol. 8365. SPIE, 2012.
  - [119] ZHU, X., AND GOLDBERG, A. *Introduction to Semi-Supervised Learning*. Morgan & Claypool, 2009.