

The Role of Principal Angles in Subspace Classification

Jiaji Huang, *Student Member, IEEE*, Qiang Qiu, Robert Calderbank, *Fellow, IEEE*

Abstract—Subspace models play an important role in a wide range of signal processing tasks, and this paper explores how the pairwise geometry of subspaces influences the probability of misclassification. When the mismatch between the signal and the model is vanishingly small, the probability of misclassification is determined by the product of the sines of the principal angles between subspaces. When the mismatch is more significant, the probability of misclassification is determined by the sum of the squares of the sines of the principal angles. Reliability of classification is derived in terms of the distribution of signal energy across principal vectors. Larger principal angles lead to smaller classification error, motivating a linear transform that optimizes principal angles. The transform presented here (TRAIT) preserves some specific characteristic of each individual class, and this approach is shown to be complementary to a previously developed transform (LRT) that enlarges inter-class distance while suppressing intra-class dispersion. Theoretical results are supported by demonstration of superior classification accuracy on synthetic and measured data even in the presence of significant model mismatch.

Index Terms—subspace, classification, SNR

I. INTRODUCTION

SIGNALS that are nominally high dimensional often exhibit a low dimensional geometric structure. For example, fixed-pose images of human faces are recorded using more than 1000 pixels, but can be represented by a 9-dimensional harmonic subspace [1]. Motion trajectories of a rigid body might be recorded by hundreds of sensors, but must intrinsically be represented by a 4-dimensional subspace [2]. There are many more examples where a low-dimensional subspace model captures intrinsic geometric structure, ranging from user ratings in a recommendation system [3] to signals emitted by multiple sources impinging at an antenna array [4]. The subspace geometry has assisted tasks of interest to both signal processing [5], [6] and machine learning communities [7], [8].

A Gaussian Mixture Model (GMM) measures proximity to a union of linear or affine subspaces, by imposing a low-rank structure on the covariance of each mixture component. It can be used to approximate a nonlinear manifold by fitting mixture components to local patches of the manifold [5], [9], hence providing a high fidelity representation of a wide variety of signal geometries. The simplicity of the model facilitates signal reconstruction [10]–[13], making GMMs a very attractive signal source model in compressed sensing. The value of low-rank GMMs extends to classification, where each class is modeled as a low-rank mixture component, and

classes are identified by their projections onto linear features. Optimal feature design is addressed in [8], [14].

The GMM is usually only an approximation to the truth. For example, the full spectrum associated with a face image follows a power law distribution, and when we truncate to the first 9 harmonic dimensions, the residual energy will be a source of error in classification. Even if the true model were a GMM, we can only learn an approximation to the true model from training data. The more data we see, the better is the fit of our empirical model, but some degree of mismatch is unavoidable. If we treat this mismatch as a form of noise, then we can use information theory to derive fundamental limits on the number of classes that can be discerned (see [15] for more details).

This paper explores how the pairwise geometry of subspaces influences the probability of misclassification. There are parallels with non-coherent wireless communication [16], where information is encoded as a subspace drawn from a fixed alphabet, and the function of the receiver is to distinguish the transmitted subspace. When each component is perfectly modeled as a Gaussian, the performance of the MAP classifier can be analyzed using the Chernoff Bound [17]. When fidelity is perfect, there is no mismatch, and fundamental limits on performance are determined by the rank of the intersection of the classes [15], [18].

In this paper, we further consider how best to discriminate classes, when the alignment between the GMM model and the data is only approximate. We make three main contributions in this paper:

- 1) We express the probability of pairwise misclassification in terms of the principal angles between the corresponding subspaces. This expression depends on the mismatch between the signal and the model. Interpreting this mismatch as noise, we provided analysis of the low, moderate, and high SNR regimes. This improves upon [18], in the sense that we have a more explicit expression of the “measurement gain” as proposed in [18].
- 2) We characterize the probability of misclassification for more general distributions near subspaces. This is motivated by the case where training samples per class are insufficient for a reliable estimate of covariance. In these cases, we have very little knowledge about the signal’s distribution and a MAP classifier is not good fit. The Nearest Subspace Classifier (NSC) provides an alternative and we use the NSC classifier rather than the MAP to bound the probability of misclassification.
- 3) We develop a feature extraction method, TRAIT, that effectively enlarges principal angles between different sub-

The authors are with the Department of Electrical and Computer Engineering, Duke University, Durham, NC, 27708 USA e-mail: jiaji.huang@duke.edu, qiang.qiu@duke.edu, robert.calderbank@duke.edu

spaces and preserves intra-class structure. We demonstrate superior classification accuracy on synthetic and measured data, particularly in the presence of significant model mismatch.

This paper is organized as follows. Section II presents the subspace geometry framework. Section III analyzes the Maximum a Posteriori (MAP) classifier under the GMM assumption. Section IV analyzes the performance of Nearest Subspace Classifier (NSC), which relaxes the GMM assumption. Section V proposes a feature extraction method, TRAIT, that exploits subspace geometry, and presents experimental results for both synthetic and measured datasets. Section VI provides a final summary.

A note on notations: we use bold upper case letters for matrices, e.g., \mathbf{X} , and bold lower case letters for vectors, e.g., \mathbf{x} . The transpose of a matrix \mathbf{X} is denoted by \mathbf{X}^\top . Scalars are written as plain letters, e.g., λ , K .

II. GEOMETRIC FRAMEWORK

Consider two subspaces \mathcal{X} and \mathcal{Y} of \mathbb{R}^n with dimensions ℓ and s respectively, where $\ell \leq s$. The principal angles between \mathcal{X} and \mathcal{Y} , denoted as $\theta_1, \dots, \theta_\ell$, are defined recursively as follows

$$\begin{aligned} \theta_1 &= \min_{\mathbf{x}_1 \in \mathcal{X}, \mathbf{y}_1 \in \mathcal{Y}} \arccos \left(\frac{\mathbf{x}_1^\top \mathbf{y}_1}{\|\mathbf{x}_1\| \|\mathbf{y}_1\|} \right), \\ &\vdots \\ \theta_j &= \min_{\substack{\mathbf{x}_j \in \mathcal{X}, \mathbf{y}_j \in \mathcal{Y} \\ \mathbf{x}_j \perp \mathbf{x}_1, \dots, \mathbf{x}_{j-1} \\ \mathbf{y}_j \perp \mathbf{y}_1, \dots, \mathbf{y}_{j-1}}} \arccos \left(\frac{\mathbf{x}_j^\top \mathbf{y}_j}{\|\mathbf{x}_j\| \|\mathbf{y}_j\|} \right), \quad j = 2, \dots, \ell. \end{aligned}$$

The vectors $\mathbf{x}_1, \dots, \mathbf{x}_\ell$ and $\mathbf{y}_1, \dots, \mathbf{y}_\ell$, are called principal vectors. The dimension of $\mathcal{X} \cap \mathcal{Y}$ is the multiplicity of zero as a principal angle. It is straightforward to compute the principal angles by calculating the singular values of $\mathbf{X}^\top \mathbf{Y}$, where \mathbf{X} and \mathbf{Y} are orthonormal bases for \mathcal{X} and \mathcal{Y} respectively. The singular values of $\mathbf{X}^\top \mathbf{Y}$ are then $\cos \theta_1, \dots, \cos \theta_\ell$.

Let $\ell = s$. The principal angles induce several distance metrics on the Grassmann manifold, of which the most widely used is the (squared) chordal distance $\mathcal{D}_c^2(\mathcal{X}, \mathcal{Y})$ [19], given by

$$\mathcal{D}_c^2(\mathcal{X}, \mathcal{Y}) = \sum_{i=1}^s \sin^2 \theta_i.$$

The chordal distance is an aggregate, and in the following sections we will see how probability of misclassification depends, not so much on this aggregate, but on the individual principal angles.

III. THE MAP CLASSIFIER FOR A GMM

We begin by considering the MAP classifier, which is optimal when the signal distribution is known. We focus on binary classification, where the two classes are equiprobable, since the generalization from two to many classes is well understood [18], [20].

We model each class as zero mean Gaussian distributed, where the covariance is near low-rank. Classification can be formulated as the following binary hypothesis testing problem

$$\begin{aligned} H_1 : \mathbf{x} &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_1) \\ H_2 : \mathbf{x} &\sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_2). \end{aligned} \quad (1)$$

We justify the zero-mean assumption by observing that in applications such as face recognition [21], or motion trajectory segmentation [2], the actual mean is considered as a nuisance parameter, and is removed prior to processing. Given the near-subspace assumption, we model the two covariances as

$$\begin{aligned} \boldsymbol{\Sigma}_1 &= \mathbf{U}_1 \boldsymbol{\Lambda}_1 \mathbf{U}_1^\top + \sigma^2 \mathbf{I} \\ \boldsymbol{\Sigma}_2 &= \mathbf{U}_2 \boldsymbol{\Lambda}_2 \mathbf{U}_2^\top + \sigma^2 \mathbf{I}. \end{aligned} \quad (2)$$

where $\mathbf{U}_1, \mathbf{U}_2 \in \mathbb{R}^{n \times d}$ are the orthonormal bases for the two signal subspaces, denoted by \mathcal{X}_1 and \mathcal{X}_2 . Typically $n \gg d$. $\boldsymbol{\Lambda}_1, \boldsymbol{\Lambda}_2 \in \mathbb{R}^{d \times d}$ are diagonal matrices of eigenvalues. We assume that the two subspaces have the same dimension d , and that the diagonal elements of $\boldsymbol{\Lambda}_1, \boldsymbol{\Lambda}_2$ are arranged in descending order. In the application to motion trajectories we take $d = 4$, and in the application to face recognition we might take $d = 9$. Denote the i -th largest eigenvalue of $\boldsymbol{\Lambda}_j$ by $\lambda_{j,i}$. Finally let σ^2 be the variance of the noise, which quantifies the degree of mismatch between the subspace model and the data.

Denote the probability of mistaking hypothesis 2 for hypothesis 1 by $\Pr(H_2|H_1)$, and define $\Pr(H_1|H_2)$ similarly. Under the assumption that the two hypotheses are equiprobable, the error probability P_e of a MAP (optimal) classifier is

$$\begin{aligned} P_e &= \frac{1}{2} [\Pr(H_2|H_1) + \Pr(H_1|H_2)] \\ &= \frac{1}{2} \int \min(\Pr(\mathbf{x}|H_1), \Pr(\mathbf{x}|H_2)) d\mathbf{x} \end{aligned} \quad (3)$$

Since this integral does not admit a closed form solution, we study the Bhattacharyya upper bound [22] to P_e instead. This bound is a special case of the Chernoff bound [17] derived using the observation $\min(a, b) \leq \sqrt{ab}$. The Bhattacharyya bound gives

$$P_e \leq \frac{1}{2} e^{-K}, \quad \text{where } K = \frac{1}{2} \ln \frac{\det(\frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2})}{\sqrt{\det \boldsymbol{\Sigma}_1 \cdot \det \boldsymbol{\Sigma}_2}}. \quad (4)$$

The numerator inside the logarithm measures the volume of space occupied by both subspaces together, and the denominator measures the volumes occupied separately. These quantities depend on the principal angles, and we now study the performance of the Bhattacharyya bound in the high, low and moderate SNR regimes.

A. The High SNR Regime

We first consider the case when $\sigma^2 \rightarrow 0$, which means that the mismatch between the signal and the model becomes vanishingly small. Since the intersection $\mathcal{X}_1 \cap \mathcal{X}_2$ between the two subspaces plays a special role, we write the two covariances as

$$\begin{aligned} \boldsymbol{\Sigma}_1 &= \mathbf{U}_{1,\cap} \boldsymbol{\Lambda}_{1,\cap} \mathbf{U}_{1,\cap}^\top + \mathbf{U}_{1,\setminus} \boldsymbol{\Lambda}_{1,\setminus} \mathbf{U}_{1,\setminus}^\top + \sigma^2 \mathbf{I}, \\ \boldsymbol{\Sigma}_2 &= \mathbf{U}_{2,\cap} \boldsymbol{\Lambda}_{2,\cap} \mathbf{U}_{2,\cap}^\top + \mathbf{U}_{2,\setminus} \boldsymbol{\Lambda}_{2,\setminus} \mathbf{U}_{2,\setminus}^\top + \sigma^2 \mathbf{I} \end{aligned} \quad (5)$$

Here both $\mathbf{U}_{1,\cap} \in \mathbb{R}^{n \times r}$ and $\mathbf{U}_{2,\cap} \in \mathbb{R}^{n \times r}$ span $\mathcal{X}_1 \cap \mathcal{X}_2$ with singular values $\boldsymbol{\Lambda}_{1,\cap}$ and $\boldsymbol{\Lambda}_{2,\cap}$ respectively. $\mathbf{U}_{1,\setminus} \in \mathbb{R}^{n \times (d-r)}$ spans $\mathcal{X}_1 \setminus \mathcal{X}_2$ with singular values $\boldsymbol{\Lambda}_{1,\setminus}$. And $\mathbf{U}_{2,\setminus} \in \mathbb{R}^{n \times (d-r)}$ spans $\mathcal{X}_2 \setminus \mathcal{X}_1$ with singular values $\boldsymbol{\Lambda}_{2,\setminus}$.

The following theorem bounds the classification error in the high SNR regime.

Theorem 1. Assume $n \geq 2(d - r)$. As $\sigma^2 \rightarrow 0$, the classification error is upper bounded as

$$P_e \leq c_1(\sigma^2)^{\frac{d-r}{2}} \left(\prod_{i=r+1}^d \sin^2 \theta_i \right)^{-\frac{1}{2}} + o\left((\sigma^2)^{\frac{d-r}{2}}\right)$$

where “ $g(\sigma^2) = o(f(\sigma^2))$ ” stands for $\lim_{\sigma^2 \rightarrow 0} \frac{g(\sigma^2)}{f(\sigma^2)} = 0$. The constant c_1 is given by,

$$c_1 = 2^{\frac{2d-r}{2}-1} \left[\frac{\text{pdet}(\mathbf{U}_{1,\cap} \mathbf{\Lambda}_{1,\cap} \mathbf{U}_{1,\cap}^\top + \mathbf{U}_{2,\cap} \mathbf{\Lambda}_{2,\cap} \mathbf{U}_{2,\cap}^\top)}{\sqrt{\prod_{i=1}^r \lambda_{1,\cap,i} \cdot \prod_{i=1}^r \lambda_{2,\cap,i}}} \cdot \prod_{i=1}^{d-r} \sqrt{\lambda_{1,\setminus,i} \cdot \lambda_{2,\setminus,i}} \right]^{-\frac{1}{2}}$$

where pdet denotes the pseudo-determinant.

Proof. The method is to expand the Bhattacharyya bound in terms of principal angles, and the details are provided in Appendix A. \square

Remark 1. 1) Typically $n \gg d$ for measured data, so the condition $n \geq 2(d - r)$ is usually satisfied.

2) The classification error is upper bounded by $(\sigma^2)^{\frac{d-r}{2}}$; the smaller the overlap between subspaces, the easier it is to discriminate between classes. When two subspaces overlap completely, there is an error floor.

There is a duality between the GMM classification problem and multiple antenna communication [23]. In multiple antenna communications, a codeword is a $d \times n$ array, where the rows are indexed by transmit antennas, the columns are indexed by time slots in a data frame, and the entries are the symbols to be transmitted. The probability of mistaking codeword C_i for codeword C_j , $\mathbb{P}\text{r}(i \rightarrow j)$, satisfies

$$\mathbb{P}\text{r}(i \rightarrow j) \leq (\sigma^2/2)^k (1/\lambda_1^2 \dots \lambda_k^2),$$

where k is the rank of $C_i - C_j$, whose singular values are $\lambda_1, \dots, \lambda_k$. The primary objective in code design for multiple antenna wireless communication is to maximize the minimum rank of the difference between distinct codewords. If the minimum rank is k , the code is said to achieve a diversity gain of k .

An important secondary objective in code design for multiple antenna wireless communication is to maximize the minimum product of the singular values of the difference between distinct codewords. This minimum product determines the coding gain.

The counterpart of coding gain in classification is the product of sines of the principal angles. This quantity determines the intercept of the error exponent with the vertical axis. The smaller the energy in the intersection of the subspaces, the smaller is the classification error. The larger the principal angles, the smaller is the classification error.

B. The Low SNR Regime

This is the case where the noise variance σ^2 and the singular values are commensurable; in other words, the mismatch between the signal and the empirical model cannot be

neglected. The MAP classifier in this case is characterized by the following theorem.

Theorem 2. When σ^2 is sufficiently large, the Bhattacharyya upper bound is sandwiched between

$$\underline{P}_e^{UB} = \frac{1}{2} \exp \left\{ -\frac{1}{\sigma^4} \left(c_2 - \frac{1}{16} \lambda_{1,1} \lambda_{2,1} \sum_{i=1}^d \cos^2 \theta_i \right) \right\}$$

and

$$\overline{P}_e^{UB} = \frac{1}{2} \exp \left\{ -\frac{1}{\sigma^4} \left(c_3 - \frac{1}{8} \lambda_{1,1} \lambda_{2,1} \sum_{i=1}^d \cos^2 \theta_i \right) \right\},$$

where $\overline{P}_e^{UB} > \underline{P}_e^{UB}$. And the constants c_2 and c_3 are given by

$$\begin{aligned} c_2 &= \frac{\sigma^4}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} - \frac{1}{2} \sum_{i=1}^d \left(\frac{\lambda_{1,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) \right] \\ &\quad + \frac{\sigma^4}{4} \left[\sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} - \frac{1}{2} \sum_{i=1}^d \left(\frac{\lambda_{2,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right] \\ c_3 &= \frac{\sigma^4}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} - \sum_{i=1}^d \left(\frac{\lambda_{1,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) \right] \\ &\quad + \frac{\sigma^4}{4} \left[\sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} - \sum_{i=1}^d \left(\frac{\lambda_{2,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right]. \end{aligned}$$

Proof. The details are given in appendix B. \square

Remark 2. The dimension of the overlap between the two subspaces plays a less important role in the low SNR regime, and classification error is a function of chordal distance. This gives rise to an interesting duality between GMM model based classification and the space-time decoding [24], where error probability is influenced by product or sum diversity in high or low SNR regime respectively.

C. The Moderate SNR Regime

We now consider a moderate noise/mismatch regime, where $\frac{p}{c(p)} \leq \frac{\lambda_{1,j}}{\sigma^2}, \frac{\lambda_{2,j}}{\sigma^2} \leq p$ for $j = 1, \dots, d$ and $p > 1, c(p) > 1$. Moderate SNR also implies that p is not very large.

The most important element in the analysis of classification error is to lower bound the term $\ln \det \left(\frac{\mathbf{\Sigma}_1 + \mathbf{\Sigma}_2}{2} \right)$ in Eq. (4),

$$\ln \det \left(\frac{\mathbf{\Sigma}_1 + \mathbf{\Sigma}_2}{2} \right) = \ln \det \left(\mathbf{I} + \frac{\mathbf{U}_1 \mathbf{\Lambda}_1 \mathbf{U}_1^\top + \mathbf{U}_2 \mathbf{\Lambda}_2 \mathbf{U}_2^\top}{2\sigma^2} \right) + n \ln \sigma^2.$$

Denote the non-zero singular values of $\mathbf{D} \triangleq \frac{1}{2\sigma^2} (\mathbf{U}_1 \mathbf{\Lambda}_1 \mathbf{U}_1^\top + \mathbf{U}_2 \mathbf{\Lambda}_2 \mathbf{U}_2^\top)$ by $\lambda_1, \dots, \lambda_{2d-r}$. Then

$$\ln \det \left(\frac{\mathbf{\Sigma}_1 + \mathbf{\Sigma}_2}{2} \right) = \sum_{i=1}^{2d-r} \ln(1 + \lambda_i) + n \ln(\sigma^2). \quad (6)$$

The following lemma provides a lower bound on $\ln(1 + \lambda_i)$.

Lemma 1. *There exists $0 \leq L < \frac{p-1}{2}$ such that for any $\lambda_i \in [L, p]$,*

$$\ln(1 + \lambda_i) \geq \ln(1 + p) + \frac{1}{1+p}(\lambda_i - p) - \frac{1}{(1+p)^2}(\lambda_i - p)^2. \quad (7)$$

Proof. See Appendix C. \square

Let $L(p)$ be the smallest possible value of L , define $c(p) = \frac{p}{2L(p)}$ if $L(p) > 0$ and $c(p) = +\infty$ if $L(p) = 0$. Note that $c(p) > 1$ since $L(p) < \frac{p-1}{2}$.

Theorem 3. *If $\frac{p}{c(p)} \leq \frac{\lambda_{1,i}}{\sigma^2}, \frac{\lambda_{2,i}}{\sigma^2} \leq p$, then the classification error is upper bounded as*

$$P_e \leq \frac{1}{2} \exp \left\{ -c_4(2d - r) + \frac{\lambda_{1,1}\lambda_{2,1}}{4\sigma^4(1+p)^2} \sum_i \cos^2 \theta_i + c_5 \right\},$$

where $c_4 = \frac{1}{2} \left[\ln(1+p) - \frac{p}{1+p} - \frac{p^2}{(1+p)^2} \right]$ and c_5 depends on p and $\frac{\lambda_{1,i}}{\sigma^2}, \frac{\lambda_{2,i}}{\sigma^2}$.

Proof. See Appendix C. \square

Remark 3. *It is straightforward to show numerically that $c(p) = 3.44, 2.79$ for $p = 4, 5$ respectively, that $c(p) \geq 2.02$ for $p \leq 10$, and that $c(p) \geq 1.61$ for $p \leq 100$. The form of the upper bound suggests that in the moderate SNR regime, the role of chordal distance is more important than the product of the sines of the principal angles.*

D. Numerical Analysis of Synthetic Data

We explore the difference between classification in the low and high SNR regimes through a simple numerical example. Consider the following pairs of subspaces:

case 1:

$$\mathbf{U}_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}^\top \quad \mathbf{U}_2 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}^\top.$$

case 2:

$$\mathbf{U}_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}^\top \quad \mathbf{U}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 0 & 0 & -1 \\ 0 & 1 & 1 & 0 \end{bmatrix}^\top.$$

We set $\mathbf{\Lambda}_1 = \mathbf{\Lambda}_2 = \mathbf{I}$ for both cases. In case 1, the two principal angles are $\theta_1 = 0, \theta_2 = \pi/2$ and in case 2, the two principal angles are $\theta_1 = \pi/4, \theta_2 = \pi/4$. The chordal distances in these two cases are the same, but in case 1 the product of sines of non-zero principal angles is 1, whereas in case 2 it is 1/2. However, there is a nontrivial intersection dimension in case 1. The product of nonzero sine principal angles is 1 for case 1, and $\frac{1}{2}$ for case 2.

We vary the degree of mismatch σ^2 , and evaluate the bounds developed in the above three theorems. In the high SNR regime, we plot the empirical misclassification probability P_e with the value $c_1(\sigma^2)^{\frac{d-r}{2}} \left(\prod_{i=r+1}^d \sin^2 \theta_i \right)^{-\frac{1}{2}}$ given in Theorem 1. In the low SNR regime, we plot the upper bound \overline{P}_e^{UB} in Theorem 2. In the moderate SNR regime, we take $p = 6$, and we vary σ^2 between $\frac{1}{p}$ and $\frac{c(p)}{p}$, so that $\frac{p}{c(p)} \leq \frac{\lambda_{1,i}}{\sigma^2}, \frac{\lambda_{2,i}}{\sigma^2} \leq p$. We then plot the upper bound in Theorem 3,

against the empirical classification error. In the high SNR regime (Fig. 1a), the classification error decays faster in Case 2 than in Case 1, consistent with Theorem 1. In the low SNR regime (Fig. 1b), there is little difference in classification error between the two cases, consistent with Theorem 2. In the moderate SNR regime (Fig. 1c), classification performance in case 1 is inferior to that in case 2, because there is a shared 1-dimensional subspace, and this is predicted by Theorem 3.

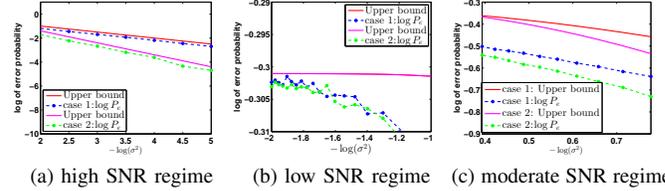


Fig. 1. Error probability as a function of the degree of mismatch. Dashed lines represent empirical estimates, and solid lines represent upper bounds. In the low SNR regime the two upper bounds coincide.

Concluding this section, we have characterized the pair-wise classification error using the principal angles between a pair of subspaces. The union bound then makes it possible to derive an upper bound on classification error for multiple classes.

IV. NEAREST SUBSPACE CLASSIFIER: EXTENDING GMM

If the class distribution is known (for example through its covariance) then the MAP classifier is optimal. If however we only know that each class is near a known low-dimensional subspace (possibly inferred from less training data) then we can substitute a Nearest Subspace Classifier (NSC) for the MAP. This Section connects performance of the NSC with principal angles, and for simplicity we focus on discriminating pairs of classes, given that the extension to multiple classes is straightforward.

Consider two classes, labeled C_1 and C_2 , distributed near two subspaces with orthonormal bases $\mathbf{U}_1, \mathbf{U}_2 \in \mathbb{R}^{n \times d}$. The NSC determines the class label of a test sample \mathbf{x} , \hat{C} , by comparing the norms of the projections onto \mathbf{U}_1 and \mathbf{U}_2 .

$$\hat{C} = \begin{cases} C_1 & \|\mathbf{U}_1^\top \mathbf{x}\|^2 \geq \|\mathbf{U}_2^\top \mathbf{x}\|^2 \\ C_2 & \text{otherwise} \end{cases}. \quad (8)$$

The preferred class label has a basis that is better aligned to the signal.

A. Derivation of the Upper Bound

Starting from the projection onto each subspace, we model the distribution of these two classes as

$$\begin{aligned} p(\mathbf{x}|C_1) &= \int p(\mathbf{x}|\boldsymbol{\alpha}, C_1) p(\boldsymbol{\alpha}) d\boldsymbol{\alpha} = \int \mathcal{N}(\mathbf{x}; \mathbf{U}_1 \boldsymbol{\alpha}, \sigma^2 \mathbf{I}) p(\boldsymbol{\alpha}) d\boldsymbol{\alpha} \\ p(\mathbf{x}|C_2) &= \int p(\mathbf{x}|\boldsymbol{\alpha}, C_2) q(\boldsymbol{\alpha}) d\boldsymbol{\alpha} = \int \mathcal{N}(\mathbf{x}; \mathbf{U}_2 \boldsymbol{\alpha}, \sigma^2 \mathbf{I}) q(\boldsymbol{\alpha}) d\boldsymbol{\alpha}. \end{aligned} \quad (9)$$

The NSC knows \mathbf{U}_1 and \mathbf{U}_2 , but is blind to $p(\boldsymbol{\alpha})$ and $q(\boldsymbol{\alpha})$, where $\boldsymbol{\alpha}$ is the expansion of the projection $\mathbf{U}_i^\top \mathbf{x}$ in the basis \mathbf{U}_i . Note that since we are not assuming a GMM, the vector $\boldsymbol{\alpha}$ need not be multivariate normal.

Let $\mathbf{V} \text{diag}\{\cos \theta_1, \dots, \cos \theta_d\} \mathbf{W}^\top$ be the singular value decomposition of $\mathbf{U}_1^\top \mathbf{U}_2$, where \mathbf{V}, \mathbf{W} are unitary, and the principal angles $\{\theta_1, \dots, \theta_d\}$ are taken in ascending order. We

may, absorb \mathbf{V} , \mathbf{W} into \mathbf{U}_1 , \mathbf{U}_2 at the cost of redefining $p(\boldsymbol{\alpha})$, $q(\boldsymbol{\alpha})$. Thus we may without loss of generality assume $\mathbf{V} = \mathbf{W} = \mathbf{I}$, i.e.,

$$\mathbf{U}_1^\top \mathbf{U}_2 = \text{diag}\{\cos \theta_1, \dots, \cos \theta_d\} \triangleq \mathbf{C}. \quad (10)$$

Define $\mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1)$ as the probability of mistaking \mathcal{C}_2 for \mathcal{C}_1 and define $\mathbb{P}\text{r}(\mathcal{C}_1|\mathcal{C}_2)$ similarly. Then the classification error is

$$P_e = \frac{1}{2}\mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1) + \frac{1}{2}\mathbb{P}\text{r}(\mathcal{C}_1|\mathcal{C}_2). \quad (11)$$

We bound $\mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1)$ using principal angles, and $\mathbb{P}\text{r}(\mathcal{C}_1|\mathcal{C}_2)$ can be analyzed in the same manner. We expand $\mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1)$ using Bayes rule as

$$\mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1) = \int \mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1, \boldsymbol{\alpha}) p(\boldsymbol{\alpha}) d\boldsymbol{\alpha}. \quad (12)$$

We bound $\mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1, \boldsymbol{\alpha})$ by writing $\mathbf{x} = \mathbf{U}_1 \boldsymbol{\alpha} + \mathbf{n}$, where the noise $\mathbf{n} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$.

$$\begin{aligned} \mathbb{P}\text{r}(\mathcal{C}_2|\mathcal{C}_1, \boldsymbol{\alpha}) &= \mathbb{P}\text{r}(\|\mathbf{U}_1^\top (\mathbf{U}_1 \boldsymbol{\alpha} + \mathbf{n})\|^2 \leq \|\mathbf{U}_2^\top (\mathbf{U}_1 \boldsymbol{\alpha} + \mathbf{n})\|^2) \\ &= \mathbb{P}\text{r}(\|\boldsymbol{\alpha} + \mathbf{U}_1^\top \mathbf{n}\|^2 \leq \|\mathbf{C}\boldsymbol{\alpha} + \mathbf{U}_2^\top \mathbf{n}\|^2), \end{aligned} \quad (13)$$

where the probability is taken w.r.t. \mathbf{n} . Denote the i -th column in \mathbf{U}_1 (\mathbf{U}_2) as $\mathbf{u}_{1,i}$ ($\mathbf{u}_{2,i}$), and the i -th element of $\boldsymbol{\alpha}$ as α_i . It follows from Eq. (13) that

$$\begin{aligned} &\mathbb{P}\text{r}(\|\boldsymbol{\alpha} + \mathbf{U}_1^\top \mathbf{n}\|^2 \leq \|\mathbf{C}\boldsymbol{\alpha} + \mathbf{U}_2^\top \mathbf{n}\|^2) \\ &= \mathbb{P}\text{r}\left(\sum_i (\alpha_i + \mathbf{u}_{1,i}^\top \mathbf{n})^2 \leq \sum_i (\cos \theta_i \alpha_i + \mathbf{u}_{2,i}^\top \mathbf{n})^2\right). \end{aligned} \quad (14)$$

We now define $a_i \triangleq \alpha_i + \mathbf{u}_{1,i}^\top \mathbf{n}$ and $b_i \triangleq \cos \theta_i \alpha_i + \mathbf{u}_{2,i}^\top \mathbf{n}$. Then Eq. (14) simplifies to

$$\begin{aligned} &\mathbb{P}\text{r}\left(\sum_i (\alpha_i + \mathbf{u}_{1,i}^\top \mathbf{n})^2 \leq \sum_i (\cos \theta_i \alpha_i + \mathbf{u}_{2,i}^\top \mathbf{n})^2\right) \\ &= \mathbb{P}\text{r}\left(\sum_i (a_i + b_i)(a_i - b_i) \leq 0\right). \end{aligned} \quad (15)$$

Lemma 2. Let a_i , b_i as defined as above. For any pair of i, j where $i \neq j$:

- 1) a_i is independent from a_j
- 2) b_i is independent from b_j
- 3) a_i is independent from b_j
- 4) $a_i + b_i$ is independent from $a_i - b_i$

Proof. The proof is given in appendix D. \square

It follows from Lemma 2 that $\sum_i (a_i + b_i)(a_i - b_i)$ is the sum of products of independently distributed normal random variables. However the product of independently distributed normal random variables need not be normal, and so we need to show that $(a_i + b_i)(a_i - b_i)$ is normally distributed.

Lemma 3 (product of normal random variable [25]). Let $x \sim \mathcal{N}(\mu_x, \sigma_x^2)$ and $y \sim \mathcal{N}(\mu_y, \sigma_y^2)$ be two independent normal variables. If $\mu_x/\sigma_x \rightarrow \infty$ and $\mu_y/\sigma_y \rightarrow \infty$ in any manner, then the distribution of xy approaches normality with mean $\mu_x \mu_y$ and variance $\mu_x^2 \sigma_y^2 + \mu_y^2 \sigma_x^2 + \sigma_x^2 \sigma_y^2$.

Applying Lemma 3 and combining the independence stated in Lemma 2, we have

Lemma 4. As $\sigma \rightarrow 0$, $\sum_i (a_i + b_i)(a_i - b_i) \sim \mathcal{N}(\sum_i \sin^2 \theta_i \alpha_i^2, 4\sigma^2 \sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2))$

Proof. The proof is given in appendix D. \square

It follows that $\mathbb{P}\text{r}(\sum_i (a_i + b_i)(a_i - b_i) \leq 0)$ is the tail probability of a normal distribution. Applying the standard tail bound, we arrive at the following theorem.

Theorem 4. As $\sigma^2 \rightarrow 0$, the classification error is upper bounded as

$$P_e \leq \int \mathcal{E}(\boldsymbol{\theta}, \boldsymbol{\alpha}, \sigma^2) \frac{p(\boldsymbol{\alpha}) + q(\boldsymbol{\alpha})}{2} d\boldsymbol{\alpha}$$

where $\mathcal{E}(\boldsymbol{\theta}, \boldsymbol{\alpha}, \sigma^2) = \frac{1}{2} \exp\left[-\frac{(\sum_{i=1}^d \sin^2 \theta_i \alpha_i^2)^2}{8\sigma^2 \sum_{i=1}^d \sin^2 \theta_i (\alpha_i^2 + \sigma^2)}\right]$.

Proof. The proof is given in appendix D. \square

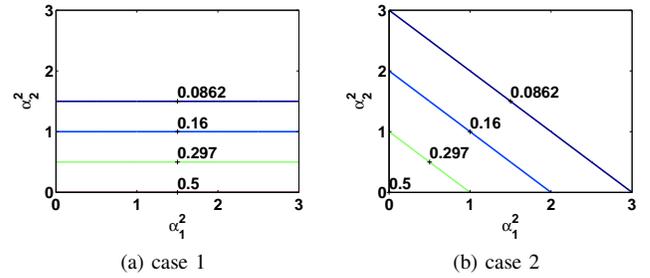


Fig. 2. Lines on which \mathcal{E} is constant for the two case studies introduced in section III-D.

We return to the two case studies introduced in Section III-D to provide some intuition about the kernel \mathcal{E} . The principal angles are $[0, \pi/2]$ in Case 1, and $[\pi/4, \pi/4]$ in Case 2. In Case 1, the kernel is constant on horizontal lines, and in Case 2, it is constant on lines of slope -1. These two cases are shown in Fig. 2, and we now make a number of general observations.

Remark 4. 1. $\mathcal{E}(\boldsymbol{\theta}, \boldsymbol{\alpha}, \sigma^2)$ is monotonically decreasing w.r.t. $\sum_i \sin^2 \theta_i \alpha_i^2$, and monotonically increasing w.r.t. σ^2 . Therefore, bigger principal angles or signal energy results in smaller classification error. Bigger noise results in bigger classification error. 2. Ignoring the higher order term of σ^2 in the denominator inside the $\exp(\cdot)$, we have

$$\mathcal{E}(\boldsymbol{\theta}, \boldsymbol{\alpha}, \sigma^2) \approx \frac{1}{2} \exp\left(-\frac{\sum_i \sin^2 \theta_i \alpha_i^2}{8\sigma^2}\right)$$

which clearly indicates that classification performance is a function of discernibility (the sine principal angles) weighted by signal energy (the α_i^2 's). 3. For fixed energy, classification error is decreased by allocating larger α_i^2 to larger θ_i .

B. Numerical Analysis of Synthetic Data

We now examine the agreement between empirical error and the upper bound given in Theorem 3. Set $n = 6$, $d = 2$,

$$\mathbf{U}_1 = [\mathbf{I}_2, \mathbf{0}_4]^\top, \quad \mathbf{U}_2 = \begin{bmatrix} \cos \theta & 0 & 0 & 0 & \sin \theta & 0 \\ 0 & \cos \theta & 0 & 0 & 0 & \sin \theta \end{bmatrix}^\top,$$

so that the two principal angles between \mathbf{U}_1 and \mathbf{U}_2 are $\theta_1 = \theta_2 = \theta$. Set $p(\boldsymbol{\alpha}) = q(\boldsymbol{\alpha}) = \mathcal{N}(\boldsymbol{\alpha}; 0, \mathbf{I}_2)$, and vary σ^2 in $[0.01, 0.5]$. Fig. 3a considers three values of θ ($\pi/6, \pi/4$, and

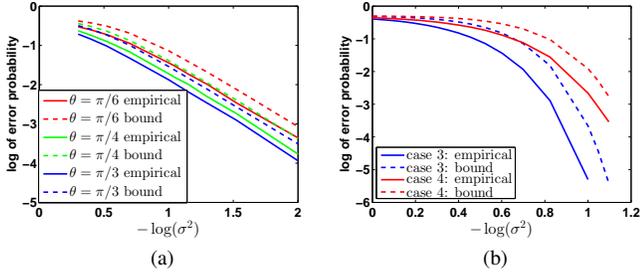


Fig. 3. Comparison of empirical NSC classification error with the upper bound obtained by numerical integration. (a) Larger principal angles reduce classification error; (b) Disproportionate assignment of signal energy to larger principal angles reduces classification error.

$\pi/3$), and shows that empirical NSC classification error tracks the upper bound obtained by numerical integration.

Next we examine the dependence of classification error on distribution of signal energy across the two modes. Set $n = 6$, $d = 2$, $\mathbf{U}_1 = [\mathbf{I}_2, \mathbf{0}_4]^\top$ and

$$\mathbf{U}_2 = \begin{bmatrix} \cos(\pi/6) & 0 & 0 & 0 & \sin(\pi/6) & 0 \\ 0 & \sin(\pi/6) & 0 & 0 & 0 & \cos(\pi/6) \end{bmatrix}^\top,$$

so that the two principal angles are $\theta_1 = \pi/6$ and $\theta_2 = \pi/3$. Fix $\|\alpha\|^2 = 1$, and compare the case when α is distributed such that $|\alpha_1| < |\alpha_2|$ (Case 3 in Fig. 3b), with the case when α is distributed such that $|\alpha_1| > |\alpha_2|$ (Case 4 in Fig. 3b). Empirical error is calculated for a range of noise variances, by randomly drawing 10,000 sample per class. Empirical NSC classification error tracks the upper bound given by numerical integration, with performance of Case 3 superior to that of Case 4.

V. TRAIT: TUNABLE RECOGNITION ADAPTED TO INTRA-CLASS TARGET

In the previous theorems, it is the principal angles that determine the performance of the classifiers in different SNR regimes. This suggests that we might improve classification by applying a linear transformation that optimizes principal angles, even at the cost of reducing dimensionality.

We denote the collection of all labeled training samples as $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_K] \in \mathbb{R}^{n \times N}$, where columns in the submatrix $\mathbf{X}_k \in \mathbb{R}^{n \times N_k}$ are samples from the k -th class. The signal subspace of \mathbf{X}_k is spanned by the orthonormal basis \mathbf{U}_k defined above. The linear transform $\mathbf{A} \in \mathbb{R}^{m \times n}$ ($m \leq n$) is designed to maximize separation of the subspaces $\mathbf{A}\mathbf{U}_1, \dots, \mathbf{A}\mathbf{U}_K$. The maximal separation is achieved when $(\mathbf{A}\mathbf{U}_j)^\top (\mathbf{A}\mathbf{U}_k) = 0$ for all $j \neq k$. In this case, all the principal angles are $\pi/2$. One approach is to use the SVD to compute the \mathbf{U}_k and then to learn the linear transformation \mathbf{A} . However we may avoid pre-computing the \mathbf{U}_k by simply encouraging $(\mathbf{A}\mathbf{X}_j)^\top (\mathbf{A}\mathbf{X}_k) = 0$ for all $j \neq k$.

We shall require that the transform \mathbf{A} preserve some specific characteristic or trait of each individual class. For example, we may target $(\mathbf{A}\mathbf{X}_k)^\top (\mathbf{A}\mathbf{X}_k) = \mathbf{X}_k^\top \mathbf{X}_k$ for all k , so that the original intra-class data structure (with noise) is preserved. Given access to a denoised signal, $\tilde{\mathbf{X}}_k$, we might instead target $(\mathbf{A}\mathbf{X}_k)^\top (\mathbf{A}\mathbf{X}_k) = \tilde{\mathbf{X}}_k^\top \tilde{\mathbf{X}}_k$ again for all k . In this case, the

intra-class dispersion due to noise is suppressed. Thus, the Gram matrix \mathbf{T} of the transformed signal can be designed to target preservation of particular intra-class structure. We formulate the optimization problem as

$$\min_{\mathbf{A} \in \mathbb{R}^{m \times n}} \frac{1}{N^2} \|(\mathbf{A}\mathbf{X})^\top (\mathbf{A}\mathbf{X}) - \mathbf{T}\|_F^2. \quad (16)$$

The block diagonal structure of the target Gram matrix \mathbf{T} promotes larger principal angles between subspaces. At the same time the diagonal blocks can be tuned to different characteristics of individual classes. For example, when side information is available, we may consider incorporating it in diagonal blocks. Here we only consider

$$\mathbf{T} = \text{diag}\{\mathbf{X}_1^\top \mathbf{X}_1, \dots, \mathbf{X}_K^\top \mathbf{X}_K\}, \quad (17)$$

as a proof-of-concept. We refer to this approach as the TRAIT algorithm, where the acronym denotes **T**unable **R**ecognition **A**dapted to **I**ntra-class **T**argets.

It is possible to minimize the objective in E.q. (16) by first minimizing $\|\mathbf{X}^\top \mathbf{P}\mathbf{X} - \mathbf{T}\|_F^2$ for $\mathbf{P} \succeq 0$ (as Proposition 1), and then factoring \mathbf{P} as $\mathbf{P} = \mathbf{A}^\top \mathbf{A}$ where $\mathbf{A} \in \mathbb{R}^{m \times n}$.

Proposition 1. *The minimizer of $\|\mathbf{X}^\top \mathbf{P}\mathbf{X} - \mathbf{T}\|_F^2$ where $\mathbf{P} \succeq 0$, is $\mathbf{P}^* = (\mathbf{X}\mathbf{X}^\top)^{-1} \mathbf{X}\mathbf{T}\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1}$.*

Proof. Proof is detailed in appendix E. \square

However when $m < n$, such a rank- m decomposition may not exist since this \mathbf{P} is not guaranteed to be rank deficient. An alternative is to learn a rank deficient \mathbf{P} by solving

$$\min_{\mathbf{P} \succeq 0} \|\mathbf{X}^\top \mathbf{P}\mathbf{X} - \mathbf{T}\|_F^2 + \lambda \|\mathbf{P}\|_*,$$

where the nuclear norm $\|\mathbf{P}\|_*$ regularizes the rank of \mathbf{P} . However this approach requires careful tuning of λ , and it is computationally more complex since we work with a matrix \mathbf{P} larger than \mathbf{A} . Given these considerations, we choose to solve (16) using gradient descent as described in Algorithm 1.

Algorithm 1 TRAIT for feature extraction

Input: labeled training samples $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_K]$, target dimension m , ($m \leq n$), target Gram matrix \mathbf{T} .

Output: feature extraction matrix (transform) $\mathbf{A} \in \mathbb{R}^{m \times n}$.

- 1: Initialize $\mathbf{A} = [\mathbf{e}_1, \dots, \mathbf{e}_m]^\top$, where \mathbf{e}_i is the i -th standard basis.
- 2: **while** stopping criteria not met **do**
- 3: Compute gradient

$$\mathbf{G} = \mathbf{A}(\mathbf{X}\mathbf{X}^\top \mathbf{A}^\top \mathbf{A}\mathbf{X}\mathbf{X}^\top - \mathbf{X}\mathbf{T}\mathbf{X}^\top).$$

- 4: Choose a positive step-size η and take a gradient step

$$\mathbf{A} \leftarrow \mathbf{A} - \eta \mathbf{G}.$$

- 5: **end while**
-

A. Related Methods

Linear Discriminant Analysis (LDA) is a classical feature extraction method which assumes each class to be Gaussian

distributed. It achieves better performance on face recognition tasks than does PCA [26]. LDA does not assume near low-rank structure of the covariances, and therefore considers a different data geometry than the one here studied.

Methods of feature extraction based on random projection have recently been developed and successfully applied to face recognition [27]. Random projection is designed to preserve pairwise distances between all data points uniformly across class labels [28].

More recently, the Low-Rank Transform (LRT) has been proposed as a method of extracting features [8]. It enlarges inter-class distance while suppressing intra-class dispersion. LRT uses the nuclear norm, $\|\mathbf{A}\mathbf{X}_i\|_*$, to measure the dispersion of the (transformed) data. The transform \mathbf{A} is

$$\arg \min_{\mathbf{A} \in \mathbb{R}^{m \times n}: \|\mathbf{A}\|_2 \leq c} \sum_{i=1}^K \|\mathbf{A}\mathbf{X}_i\|_* - \|\mathbf{A}\mathbf{X}\|_*.$$

What motivates the choice of the nuclear norm is that it is the convex relaxation of rank [8]. In the high SNR regime, Theorem 1 suggests that classification error decreases when the union of subspaces has large rank. LRT encourages the rank of the union to be large, and it works well in a regime where model mismatch is small. Experiments presented in Section V-C suggest that TRAIT may be more robust to model mismatch (Fig. 8).

B. Two Properties of the TRAIT Transform

On synthetic and measured data, we show that TRAIT effectively enlarges the angles between different subspaces and preserves intra-class structure. We also compare the classification accuracy of features extracted by TRAIT and the methods in Section V-A. For synthetic data, the class distribution is known exactly, and the MAP classifier is used to measure classification accuracy. For measured data, the class distribution is unknown a priori, and the NSC classifier is employed instead.

1) *Enlargement of the Principal angles:* The synthetic dataset has parameters $n = 10$, $d = 1$ and $K = 3$.

$$\Sigma_k = \mathbf{U}_k \mathbf{U}_k^\top + 10^{-2} \mathbf{I} (k = 1, 2, 3),$$

where \mathbf{U}_k is a normalized n -vector with i.i.d. Gaussian random entries. Samples of the k -th class are i.i.d drawn from $\mathcal{N}(\mathbf{0}, \Sigma_k)$. For each class, 100 samples are used for learning the transform and 10000 are used for testing. On the training data, we learn the transform respectively via LDA, LRT, and TRAIT with target dimension $m = 3, \dots, 10$. Then on each test datum, we apply the learned transforms as well as random projection (each entry drawn from $\mathcal{N}(0, 1)$) and classify using a MAP classifier.

We visualize original and transformed data via projection (PCA basis) into 3-dimensional Euclidean space. When the target feature dimension $m = 3$, the results are shown in Fig. 4. Each class is represented by a different color. After transforming the data, we use the SVD to calculate the basis vector ($d = 1$) that best describes each class, and we calculate the pairwise angles between basis vectors. The pairwise angles

are significantly increased by both LRT and TRAIT. By contrast, neither LDA nor random projection increase separation between one-dimensional subspaces.

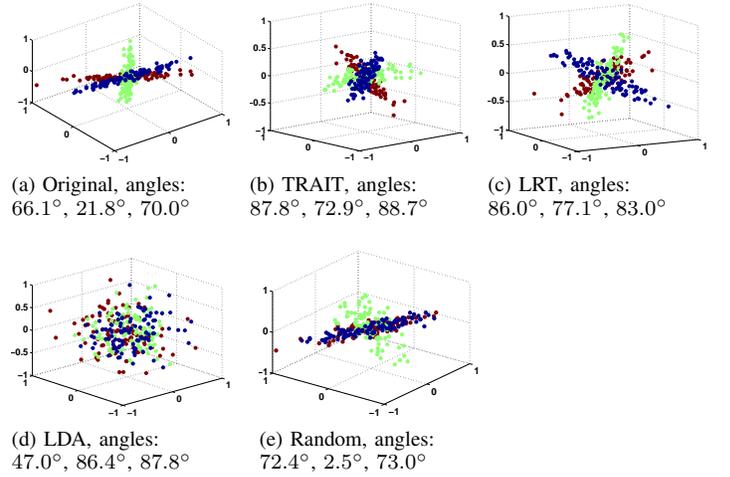


Fig. 4. Embeddings of original and transformed data.

We now vary the feature dimension m , and compare the error probability of the MAP classifier across the different methods of extracting features. Fig. 5 shows that the performance of TRAIT and LRT are similar, and that both are superior to LDA and random projection. Note that after dimension reduction, TRAIT is still able to match error probabilities achieved with the original data.

2) *Preservation of Intra-class Structure:* When a convex body, e.g., human face, is illuminated, the resulting image is represented by spherical harmonics. It has been shown that a 9-dimensional subspace is sufficient to capture the geometry of an individual subject [1]. The extended Yale B face database includes 38 subjects, each with 64 images taken under different illumination conditions. We use a cropped version of this data set¹, where each image is of size $32 \times 32 = 1024$.

For each subject, we randomly select half of the 64 images for training, and retain the other half for testing. For all feature extraction methods, we vary the target dimension m , and apply the NSC to the transformed data. The NSC achieves much higher accuracy on features extracted by TRAIT and LRT (Fig. 6).

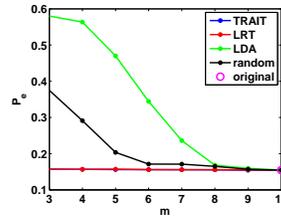


Fig. 5. MAP classifier's P_e on synthetic data. Note that TRAIT (blue) and LRT (red) almost overlap.

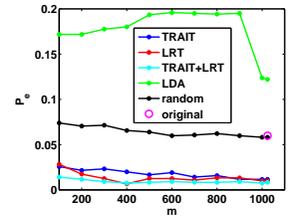


Fig. 6. NSC's P_e on original/transformed face images. Concatenation of TRAIT and LRT features (TRAIT+LRT) provides superior results

¹<http://www.cad.zju.edu.cn/home/dengcai/Data/FaceData.html>

We also observe in Fig. 7 that the features extracted by TRAIT and LRT are quite different, suggesting that information present in one view is somewhat independent of information present in the other. This is confirmed by applying NSC to the concatenation of the two views (TRAIT+LRT in Fig. 6), and observing that classification accuracy is increased.

The intra-class structure preserving property of TRAIT is evident in Fig. 7 where we view transformed classes as faces in the original image domain. The original images of subject 10 are displayed together with their TRAIT and LRT transforms. TRAIT preserves a diversity of illumination conditions, whereas LRT blurs the differences between images. Classification performance is improved by using LRT and TRAIT features in combination.

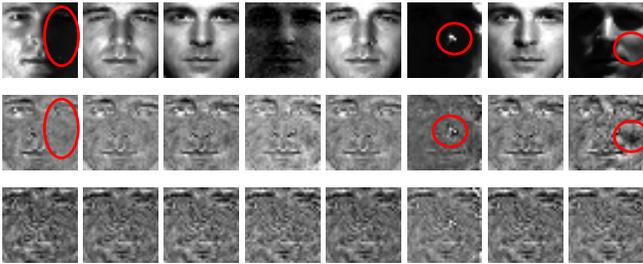


Fig. 7. Comparison of original images (top) with TRAIT transformed images (middle) and LRT transformed images (bottom). Red circles indicate structure that is present in both the original and the TRAIT transformed image.

C. Robustness to Model Mismatch

In the previous sections, we have demonstrated the effectiveness of TRAIT and LRT on both synthetic and real data. In this section, we present experiments showing that TRAIT is more robust with respect to model mismatch than is LRT. In many real world problems, data may not be exactly GMM distributed. Even if they are, there may not be sufficient training data to learn the covariances. Therefore, we use NSC throughout this section to assess the discriminability of the extracted features. Moreover, having seen the effectiveness of dimension reduction in previous sections, we turn to learning dimension reduced features, thereby saving computational cost on measured datasets.

1) *Synthetic Data*: The synthetic data is a three-class dataset, where datum $\mathbf{x} \in \mathbb{R}^{100}$ in the k -th ($k = 1, 2, 3$) class is generated as

$$\mathbf{x} = \mathbf{U}_k \boldsymbol{\alpha} + \mathbf{n},$$

with $\mathbf{U}_k \in \mathbb{R}^{100 \times 5}$ and $\mathbf{U}_k^T \mathbf{U}_k = \mathbf{I}$. $\boldsymbol{\alpha} \sim \text{Uniform}[-2, 2]$ and $\mathbf{n} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{100})$. Note the data is not GMM distributed. Each class has 100 training samples and 10000 testing samples. We vary σ^2 and use NSC to classify TRAIT and LRT extracted features. Here we fix the extracted feature dimension to be 30.

Fig. 8 shows the NSC classification accuracy as a function of σ . Both TRAIT and LRT significantly improves classification performance compared with no transform. However, with increasing noise, TRAIT features outperform LRT features, showing greater robustness to model mismatch.

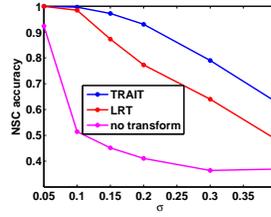


Fig. 8. NSC performance on TRAIT and LRT features under different SNR



Fig. 9. From top to bottom row: subjects in PIE, UMIST and ORL database, taken under different poses

2) *Face Images with non-frontal Poses*: It is known that human frontal face images are well modeled by subspaces. For example, the Yale-B face in section V-B2, where LRT slightly outperforms TRAIT. Now we further compare the performance of TRAIT and LRT in more mismatched cases by introducing non-frontal face images. We validate performance on three publicly available datasets, PIE [29], UMIST² and ORL³. All of them have a considerable number of non-frontal face images. Fig. 9 shows one subject from each database with different poses.

The PIE dataset includes 18562 64×48 images of 68 subjects. Each image is labeled with one of 13 different pose tags. We randomly select 7 pose tags and the images of these tags are used as training samples. The rest are used in testing. UMIST comprises 575 112×92 images of 20 subjects, and ORL comprises 400 112×92 images of 40 subjects. These two datasets have no pose tags. We split the UMIST and ORL datasets using the strategy followed for the Yale-B dataset in Section V-B2. We derive 1000-dimensional features for each of random projection, LDA, LRT and TRAIT. Table I lists accuracies of NSC classification for the different algorithms.

TABLE I
NSC ACCURACY ON ORIGINAL AND 1000 DIMENSIONAL (COMPRESSED) EXTRACTED FEATURES

	PIE	UMIST	ORL
Original	74.57%	96.14%	95.50%
random	72.14%	95.44%	94.50%
LDA	40.10%	84.91%	92.00%
LRT	70.80%	96.84%	95.00%
TRAIT	76.11%	97.90%	97.00%

In all cases, TRAIT has the highest classification accuracy and outperforms LRT. LRT optimizes the rank (its convex relaxation), which is critical for reducing classification error in the high SNR regime. However, in this low SNR regime, TRAIT gains more discrimination via explicitly “orthogonalizing” between the classes. The criteria employed by TRAIT do not depend on the specific SNR regime and therefore are more robust.

VI. CONCLUSION

In a low-rank Gaussian Mixture Model, we have explored how the probability of misclassification is governed by prin-

²<http://www.sheffield.ac.uk/eee/research/iel/research/face>

³<http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>

principal angles between subspaces. In the low-noise regime, the Bhattacharyya upper bound on misclassification is determined by the product of the sines of the principal angles. In the high/moderate-noise regime it is determined by the sum of the squares of the sines of the principal angles. Analysis of the Nearest Subspace Classifier connected reliability of classification to the distribution of signal energy across principal vectors. Classification was shown to be more reliable when more signal energy is associated with principal vectors corresponding to large principal angles. This observation motivated the design of a transform, TRAIT, that achieves superior classification performance by enlarging principal angles and preserving intra-class structure. Finally we showed that TRAIT complements a prior approach that enlarge inter-class distance while suppressing intra-class dispersion, and that it is more robust to model mismatch.

APPENDIX A PROOF OF HIGH SNR CASE

Proof of Theorem 1 We have

$$\begin{aligned}\det \mathbf{\Sigma}_1 &= (\sigma^2)^{n-d} \prod_{i=1}^d (\lambda_{1,i} + \sigma^2), \\ \det \mathbf{\Sigma}_2 &= (\sigma^2)^{n-d} \prod_{i=1}^d (\lambda_{2,i} + \sigma^2).\end{aligned}$$

Let the SVD of $\mathbf{U}_{1,\cap} \mathbf{\Lambda}_{1,\cap} \mathbf{U}_{1,\cap}^\top + \mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus} \mathbf{U}_{1,\setminus}^\top + \mathbf{U}_{2,\cap} \mathbf{\Lambda}_{2,\cap} \mathbf{U}_{2,\cap}^\top + \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus} \mathbf{U}_{2,\setminus}^\top$ be $\mathbf{Z} \mathbf{\Lambda} \mathbf{Z}^\top$, where $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_{2d-r}\}$. Then,

$$\det \left(\frac{\mathbf{\Sigma}_1 + \mathbf{\Sigma}_2}{2} \right) = (\sigma^2)^{n-2d+r} \prod_{i=1}^{2d-r} \left(\frac{\lambda_i}{2} + \sigma^2 \right).$$

Substituting the above into the Bhattacharyya bound, we have

$$\begin{aligned}P_e &\leq \frac{1}{2} (\sigma^2)^{\frac{d-r}{2}} \cdot \left[\frac{\sqrt{\prod_{i=1}^d (\lambda_{1,i} + \sigma^2) \prod_{i=1}^d (\lambda_{2,i} + \sigma^2)}}{\prod_{i=1}^{2d-r} \left(\frac{\lambda_i}{2} + \sigma^2 \right)} \right]^{\frac{1}{2}} \\ &= (\sigma^2)^{\frac{d-r}{2}} \cdot 2^{\frac{2d-r}{2}-1} \left[\frac{\sqrt{\prod_{i=1}^d \lambda_{1,i} \prod_{i=1}^d \lambda_{2,i}}}{\prod_{i=1}^{2d-r} \lambda_i} \right]^{\frac{1}{2}} \\ &\quad + o \left((\sigma^2)^{\frac{d-r}{2}} \right).\end{aligned}\tag{18}$$

Our objective is to expand $\prod_{i=1}^{2d-r} \lambda_i$ in terms of principal angles. Since the image of $\mathbf{U}_{1,\cap}$ (or $\mathbf{U}_{2,\cap}$) is orthogonal to $\mathbf{U}_{1,\setminus}$ and $\mathbf{U}_{2,\setminus}$,

$$\begin{aligned}\prod_{i=1}^{2d-r} \lambda_i &= \text{pdet}(\mathbf{U}_{1,\cap} \mathbf{\Lambda}_{1,\cap} \mathbf{U}_{1,\cap}^\top + \mathbf{U}_{2,\cap} \mathbf{\Lambda}_{2,\cap} \mathbf{U}_{2,\cap}^\top) \\ &\quad \cdot \text{pdet}([\mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \quad \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}}][\mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \quad \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}}]^\top) \\ &= \text{pdet}(\mathbf{U}_{1,\cap} \mathbf{\Lambda}_{1,\cap} \mathbf{U}_{1,\cap}^\top + \mathbf{U}_{2,\cap} \mathbf{\Lambda}_{2,\cap} \mathbf{U}_{2,\cap}^\top) \\ &\quad \cdot \det([\mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \quad \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}}]^\top [\mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \quad \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}}]),\end{aligned}$$

where we assume $n \geq 2(d-r)$ in order to derive the second equality, which simplifies as follows:

$$\begin{aligned}&\det([\mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \quad \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}}]^\top [\mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \quad \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}}]) \\ &= \det \left(\begin{bmatrix} \mathbf{\Lambda}_{1,\setminus} & & & \\ & \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \mathbf{U}_{1,\setminus}^\top \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}} & & \\ & & \mathbf{\Lambda}_{2,\setminus} & \\ & & & \end{bmatrix} \right) \\ &= \det(\mathbf{\Lambda}_{1,\setminus}) \det(\mathbf{\Lambda}_{2,\setminus} - \\ &\quad \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}} \mathbf{U}_{2,\setminus}^\top \mathbf{U}_{1,\setminus} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \mathbf{\Lambda}_{1,\setminus}^{-1} \mathbf{\Lambda}_{1,\setminus}^{\frac{1}{2}} \mathbf{U}_{1,\setminus}^\top \mathbf{U}_{2,\setminus} \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}}) \\ &= \det(\mathbf{\Lambda}_{1,\setminus}) \det \left(\mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}} (\mathbf{I} - \mathbf{U}_{2,\setminus}^\top \mathbf{U}_{1,\setminus} \mathbf{U}_{1,\setminus}^\top \mathbf{U}_{2,\setminus}) \mathbf{\Lambda}_{2,\setminus}^{\frac{1}{2}} \right) \\ &= \prod_{i=1}^{d-r} \lambda_{1,\setminus,i} \cdot \prod_{i=1}^{d-r} \lambda_{2,\setminus,i} \cdot \prod_{i=r+1}^d \sin^2 \theta_i.\end{aligned}\tag{19}$$

The last equality follows from the observation that the eigenvalues of $\mathbf{U}_{2,\setminus}^\top \mathbf{U}_{1,\setminus} \mathbf{U}_{1,\setminus}^\top \mathbf{U}_{2,\setminus}$ are $\cos^2 \theta_{r+1}, \dots, \cos^2 \theta_d$. The theorem now follows by substituting Eq. (19) into Eq. (18). \square

APPENDIX B PROOF OF LOW SNR CASE

We first state and prove (for completeness) two preliminary lemmas that are needed to characterize classification error.

Lemma 5. Let $\mathbf{D} \in \mathbb{R}^{n \times n}$ be any positive semi-definite matrix with all eigenvalues smaller than 1, then

$$\text{tr}(\mathbf{D}) - \frac{1}{2} \text{tr}(\mathbf{D}^2) \leq \ln \det(\mathbf{I}_n + \mathbf{D}) \leq \text{tr}(\mathbf{D}) - \frac{1}{4} \text{tr}(\mathbf{D}^2).$$

Proof. Denote the nonnegative eigenvalues of $\mathbf{D} \succeq 0$ as d_1, \dots, d_n , where $d_1, \dots, d_n \leq 1$. Then

$$\ln \det(\mathbf{I}_n + \mathbf{D}) = \ln \prod_{i=1}^n (1 + d_i) = \sum_{i=1}^n \ln(1 + d_i).$$

Since $x - \frac{x^2}{2} \leq \ln(1 + x) \leq x - \frac{x^2}{4}$ for all $x \in [0, 1]$, we obtain

$$\sum_i d_i - \frac{d_i^2}{2} \leq \ln \det(\mathbf{I}_n + \mathbf{D}) \leq \sum_i d_i - \frac{d_i^2}{4},$$

which reduces to

$$\text{tr}(\mathbf{D}) - \frac{1}{2} \text{tr}(\mathbf{D}^2) \leq \ln \det(\mathbf{I}_n + \mathbf{D}) \leq \text{tr}(\mathbf{D}) - \frac{1}{4} \text{tr}(\mathbf{D}^2).$$

This bound is very tight when all the d_i 's approach 0. \square

Lemma 6. Suppose $\mathbf{U} \in \mathbb{R}^{n \times d}$, $\mathbf{V} \in \mathbb{R}^{n \times d}$ are two orthonormal bases and that $\mathbf{\Phi} \in \mathbb{R}^{d \times d}$, $\mathbf{\Psi} \in \mathbb{R}^{d \times d}$ are diagonal with nonnegative decreasing diagonal elements ϕ_1, \dots, ϕ_d and ψ_1, \dots, ψ_d respectively. Denote the i -th principal angle between \mathbf{U} and \mathbf{V} as θ_i where $i = 1, \dots, d$. Then

$$\phi_d \psi_d \sum_i \cos^2 \theta_i \leq \text{tr}(\mathbf{U} \mathbf{\Phi} \mathbf{U}^\top \mathbf{V} \mathbf{\Psi} \mathbf{V}^\top) \leq \phi_1 \psi_1 \sum_i \cos^2 \theta_i.$$

Proof. Let the Singular Value Decomposition of $\mathbf{U}^\top \mathbf{V}$ be $\mathbf{J} \mathbf{C} \mathbf{H}^\top$, then $\text{tr}(\mathbf{U}^\top \mathbf{V} \mathbf{V}^\top \mathbf{U}) = \text{tr}(\mathbf{C}^2) = \sum_i \cos^2 \theta_i$. We have

$$\begin{aligned}\text{tr}(\mathbf{U} \mathbf{\Phi} \mathbf{U}^\top \mathbf{V} \mathbf{\Psi} \mathbf{V}^\top) &= \text{tr}(\mathbf{\Phi} \mathbf{U}^\top \mathbf{V} \mathbf{\Psi} \mathbf{V}^\top \mathbf{U}) \\ &= \text{tr}(\mathbf{\Phi} \mathbf{J} \mathbf{C} \mathbf{H}^\top \mathbf{\Psi} \mathbf{H} \mathbf{C} \mathbf{J}^\top) = \text{tr}(\mathbf{J}^\top \mathbf{\Phi} \mathbf{J} \mathbf{C} \mathbf{H}^\top \mathbf{\Psi} \mathbf{H} \mathbf{C}).\end{aligned}$$

For any two positive semidefinite matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times m}$, let the maximum and minimum eigenvalues of \mathbf{A} be $\lambda_1(\mathbf{A}), \lambda_m(\mathbf{A})$ respectively, then by [30]

$$\lambda_m(\mathbf{A}) \operatorname{tr}(\mathbf{B}) \leq \operatorname{tr}(\mathbf{AB}) \leq \lambda_1(\mathbf{A}) \operatorname{tr}(\mathbf{B}).$$

Hence,

$$\begin{aligned} \operatorname{tr}(\mathbf{U}\Phi\mathbf{U}^\top\mathbf{V}\Psi\mathbf{V}^\top) &\leq \phi_1 \operatorname{tr}(\mathbf{CH}^\top\Psi\mathbf{HC}) = \phi_1 \operatorname{tr}(\mathbf{H}^\top\Psi\mathbf{HC}^2) \\ &\leq \phi_1\psi_1 \operatorname{tr}(\mathbf{C}^2) = \phi_1\psi_1 \sum_i \cos^2 \theta_i. \end{aligned}$$

The lower bound can be proved in the same way. This bound becomes tight when the diagonal elements of Φ and Ψ are uniform. \square

Proof of Theorem 2 We are now ready to prove theorem 2. We expand K in Eq. (4) as

$$K = \frac{1}{2} \ln \det \left(\frac{\Sigma_1 + \Sigma_2}{2} \right) - \frac{1}{4} (\ln \det \Sigma_1 + \ln \det \Sigma_2). \quad (20)$$

The second term becomes:

$$-\frac{1}{4} \left[\sum_{i=1}^d \ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) + \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right] - \frac{n}{2} \ln(\sigma^2), \quad (21)$$

and we use Lemma 5 to bound the first term. Note that

$$\begin{aligned} &\frac{1}{2} \ln \det \left(\frac{\Sigma_1 + \Sigma_2}{2} \right) \\ &= \frac{1}{2} \ln \det \left[\sigma^2 \left(\mathbf{I} + \frac{\mathbf{U}_1\Lambda_1\mathbf{U}_1^\top + \mathbf{U}_2\Lambda_2\mathbf{U}_2^\top}{2\sigma^2} \right) \right] \\ &= \frac{n}{2} \ln(\sigma^2) + \frac{1}{2} \ln \det \left(\mathbf{I} + \frac{\mathbf{U}_1\Lambda_1\mathbf{U}_1^\top + \mathbf{U}_2\Lambda_2\mathbf{U}_2^\top}{2\sigma^2} \right). \end{aligned} \quad (22)$$

Let $\mathbf{D} \triangleq \frac{\mathbf{U}_1\Lambda_1\mathbf{U}_1^\top + \mathbf{U}_2\Lambda_2\mathbf{U}_2^\top}{2\sigma^2}$. We apply Lemma 5 to bound $\frac{1}{2} \ln \det \left(\frac{\Sigma_1 + \Sigma_2}{2} \right)$:

$$\begin{aligned} &\frac{n}{2} \ln(\sigma^2) + \frac{1}{2} \left[\operatorname{tr}(\mathbf{D}) - \frac{1}{2} \operatorname{tr}(\mathbf{D}^2) \right] \leq \frac{1}{2} \ln \det \left(\frac{\Sigma_1 + \Sigma_2}{2} \right) \\ &\leq \frac{n}{2} \ln(\sigma^2) + \frac{1}{2} \left[\operatorname{tr}(\mathbf{D}) - \frac{1}{4} \operatorname{tr}(\mathbf{D}^2) \right], \end{aligned} \quad (23)$$

Expanding $\operatorname{tr}(\mathbf{D})$ gives

$$\begin{aligned} &\frac{n}{2} \ln(\sigma^2) + \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} + \sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} \right] - \frac{1}{4} \operatorname{tr}(\mathbf{D}^2) \\ &\leq \frac{1}{2} \ln \det \left(\frac{\Sigma_1 + \Sigma_2}{2} \right) \\ &\leq \frac{n}{2} \ln(\sigma^2) + \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} + \sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} \right] - \frac{1}{8} \operatorname{tr}(\mathbf{D}^2). \end{aligned} \quad (24)$$

Note that

$$\operatorname{tr}(\mathbf{D}^2) = \frac{1}{4\sigma^4} \left(\sum_{i=1}^d \lambda_{1,i}^2 + \sum_{i=1}^d \lambda_{2,i}^2 + 2 \operatorname{tr}(\mathbf{U}_1\Lambda_1\mathbf{U}_1^\top\mathbf{U}_2\Lambda_2\mathbf{U}_2^\top) \right). \quad (25)$$

Invoking Lemma 6 to bound the last term of the above:

$$\begin{aligned} \operatorname{tr}(\mathbf{U}_1\Lambda_1\mathbf{U}_1^\top\mathbf{U}_2\Lambda_2\mathbf{U}_2^\top) &\geq \lambda_{1,d}\lambda_{2,d} \sum_i \cos^2 \theta_i \\ \operatorname{tr}(\mathbf{U}_1\Lambda_1\mathbf{U}_1^\top\mathbf{U}_2\Lambda_2\mathbf{U}_2^\top) &\leq \lambda_{1,1}\lambda_{2,1} \sum_i \cos^2 \theta_i \end{aligned} \quad (26)$$

Combining Eq. (21) to (26), we obtain upper and lower bounds on K ,

1) *Upper bound:*

$$\begin{aligned} K &\leq \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} + \sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} \right] \\ &\quad - \frac{1}{32\sigma^4} \left(\sum_{i=1}^d \lambda_{1,i}^2 + \sum_{i=1}^d \lambda_{2,i}^2 + 2\lambda_{1,d}\lambda_{2,d} \sum_{i=1}^d \cos^2 \theta_i \right) \\ &\quad - \frac{1}{4} \left[\sum_{i=1}^d \ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) + \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right] \\ &= \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} - \frac{1}{2} \sum_{i=1}^d \left(\frac{\lambda_{1,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) \right] \\ &\quad + \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} - \frac{1}{2} \sum_{i=1}^d \left(\frac{\lambda_{2,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right] \\ &\quad - \frac{1}{16\sigma^4} \lambda_{1,d}\lambda_{2,d} \sum_{i=1}^d \cos^2 \theta_i \\ &\triangleq \frac{1}{\sigma^4} \left(c_2 - \frac{1}{16} \lambda_{1,d}\lambda_{2,d} \sum_{i=1}^d \cos^2 \theta_i \right). \end{aligned} \quad (27)$$

2) *Lower bound:*

$$\begin{aligned} K &\geq \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} + \sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} \right] \\ &\quad - \frac{1}{16\sigma^4} \left(\sum_{i=1}^d \lambda_{1,i}^2 + \sum_{i=1}^d \lambda_{2,i}^2 + 2\lambda_{1,1}\lambda_{2,1} \sum_{i=1}^d \cos^2 \theta_i \right) \\ &\quad - \frac{1}{4} \left[\sum_{i=1}^d \ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) + \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right] \\ &= \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{1,i}}{\sigma^2} - \sum_{i=1}^d \left(\frac{\lambda_{1,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) \right] \\ &\quad + \frac{1}{4} \left[\sum_{i=1}^d \frac{\lambda_{2,i}}{\sigma^2} - \sum_{i=1}^d \left(\frac{\lambda_{2,i}}{2\sigma^2} \right)^2 - \sum_{i=1}^d \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right] \\ &\quad - \frac{1}{8\sigma^4} \lambda_{1,1}\lambda_{2,1} \sum_{i=1}^d \cos^2 \theta_i \\ &\triangleq \frac{1}{\sigma^4} \left(c_3 - \frac{1}{8} \lambda_{1,d}\lambda_{2,d} \sum_{i=1}^d \cos^2 \theta_i \right). \end{aligned} \quad (28)$$

Negating K and exponentiating gives theorem 2. \square

APPENDIX C

PROOF OF MODERATE SNR CASE

Proof of Lemma 1 consider the function

$$f(\lambda_i) = \ln(1+\lambda_i) - \ln(1+p) - \frac{1}{1+p}(\lambda_i-p) + \frac{1}{(1+p)^2}(\lambda_i-p)^2,$$

defined in $[0, p]$. Its derivative is

$$f'(\lambda_i) = \frac{1}{1+\lambda_i} - \frac{1}{1+p} + \frac{2(\lambda_i-p)}{(1+p)^2} = \frac{(p-\lambda_i)(p-1-2\lambda_i)}{(1+\lambda_i)(1+p)^2},$$

which is positive in $[0, \frac{p-1}{2})$ and negative in $(\frac{p-1}{2}, p]$. Therefore, $f(\lambda_i)$ is monotonically increasing in $[0, \frac{p-1}{2})$ and decreasing in $(\frac{p-1}{2}, p]$. Further, $f(p) = 0$ and $f(0) = -\ln(1+p) + \frac{p}{1+p} + \frac{p^2}{(1+p)^2}$ whose sign depends on the value of p . The shape of $f(\lambda_i)$ is now characterized. There exists $L < \frac{p-1}{2}$ such that $f(\lambda_i) \geq 0$ when $\lambda_i \in [L, p]$. \square

Before proving theorem 3, we need to bound λ_i using Weyl's inequality [31].

Lemma 7 (Weyl's inequality [31]). *Let \mathbf{M} and \mathbf{P} be two $n \times n$ Hermitian matrices, with eigenvalues $\mu_1 \geq \dots \geq \mu_n$ and $\nu_1 \geq \dots \geq \nu_n$ respectively. Denote the eigenvalues of $\mathbf{M} + \mathbf{P}$ by $\gamma_1 \geq \dots \geq \gamma_n$. Then*

$$\max(\mu_i + \nu_n, \nu_i + \mu_n) \leq \gamma_i \leq \min(\mu_i + \nu_1, \nu_i + \mu_1).$$

Proof of Theorem 3 Since $\frac{p}{c(p)} \leq \frac{\lambda_{1,i}}{\sigma^2}, \frac{\lambda_{2,i}}{\sigma^2} \leq p$, by the Weyl's inequality, $\frac{p}{2c(p)} = \frac{p/c(p)+0}{2} \leq \lambda_i \leq \frac{p+p}{2} = p$. Further, since $1 \leq c(p) \leq \frac{2L(p)}{p}$, we have $\lambda_1, \dots, \lambda_{2d-r} \in [L(p), p]$. By definition of $L(p)$, we can invoke Eq. (7) in Lemma 1 to obtain

$$\begin{aligned} \ln \det \left(\frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2} \right) &= \sum_{i=1}^{2d-r} \ln(1 + \lambda_i) + n \ln(\sigma^2) \\ &\geq (2d-r) \ln(1+p) + \frac{\text{tr } \mathbf{D} - p(2d-r)}{1+p} \\ &\quad - \frac{\text{tr } \mathbf{D}^2 - 2p \text{tr } \mathbf{D} + p^2(2d-r)}{(1+p)^2} + n \ln(\sigma^2). \end{aligned} \quad (29)$$

Notice $\text{tr } \mathbf{D} = \frac{1}{2} \sum_i (\frac{\lambda_{1,i}}{\sigma^2} + \frac{\lambda_{2,i}}{\sigma^2})$, and by Eq. (25) and (26), $\text{tr } \mathbf{D}^2 \leq \frac{1}{4\sigma^4} (\sum_i \lambda_{1,i}^2 + \lambda_{2,i}^2 + 2\lambda_{1,1}\lambda_{2,1} \sum_i \cos^2 \theta_i)$. Substituting these into Eq. (29), we get

$$\begin{aligned} &\ln \det \left(\frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2} \right) \\ &\geq n \ln(\sigma^2) + (2d-r) \left[\ln(1+p) - \frac{p}{1+p} - \frac{p^2}{(1+p)^2} \right] \\ &\quad + \frac{1+3p}{2\sigma^2(1+p)^2} \left(\sum_i \lambda_{1,i} + \lambda_{2,i} \right) \\ &\quad - \frac{1}{4\sigma^4(1+p)^2} \left(\sum_i \lambda_{1,i}^2 + \lambda_{2,i}^2 + 2\lambda_{1,1}\lambda_{2,1} \sum_i \cos^2 \theta_i \right) \end{aligned}$$

Substituting the above into the Bhattacharyya bound (4) yields an upper bound on P_e , of the form given in Theorem 3. In particular,

$$c_4 = \frac{1}{2} \left[\ln(1+p) - \frac{p}{1+p} - \frac{p^2}{(1+p)^2} \right],$$

and

$$\begin{aligned} c_5 &= -\frac{1+3p}{4\sigma^2(1+p)^2} \sum_i (\lambda_{1,i} + \lambda_{2,i}) + \frac{\sum_i \lambda_{1,i}^2 + \lambda_{2,i}^2}{8\sigma^4(1+p)^2} \\ &\quad + \frac{1}{4} \sum_i \left[\ln \left(1 + \frac{\lambda_{1,i}}{\sigma^2} \right) + \ln \left(1 + \frac{\lambda_{2,i}}{\sigma^2} \right) \right]. \end{aligned}$$

\square

APPENDIX D ANALYSIS OF NSC

Proof of Lemma 2 Since that the joint distribution of $[a_i \ a_j]^\top$, $[b_i \ b_j]^\top$, $[a_i \ b_j]^\top$ and $[a_i + b_i \ a_i - b_i]^\top$ are all Gaussian, it suffices to show that all covariance are diagonal. For any $i \neq j$,

$$\begin{aligned} \begin{bmatrix} a_i \\ a_j \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} \alpha_i \\ \alpha_j \end{bmatrix}, \sigma^2 \mathbf{I}_2 \right) & \begin{bmatrix} b_i \\ b_j \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} \cos \theta_i \alpha_i \\ \cos \theta_j \alpha_j \end{bmatrix}, \sigma^2 \mathbf{I}_2 \right) \\ \begin{bmatrix} a_i \\ b_i \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} \alpha_i \\ \cos \theta_j \alpha_j \end{bmatrix}, \sigma^2 \mathbf{I}_2 \right). \end{aligned}$$

For any i ,

$$\begin{aligned} \begin{bmatrix} a_i \\ b_i \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} \alpha_i \\ \cos \theta_i \alpha_i \end{bmatrix}, \sigma^2 \begin{bmatrix} 1 & \cos \theta_i \\ \cos \theta_i & 1 \end{bmatrix} \right) \\ \begin{bmatrix} a_i + b_i \\ a_i - b_i \end{bmatrix} &\sim \mathcal{N} \left(\begin{bmatrix} (1 + \cos \theta_i) \alpha_i \\ (1 - \cos \theta_i) \alpha_i \end{bmatrix}, \right. \\ &\quad \left. 2\sigma^2 \begin{bmatrix} (1 + \cos \theta_i) & 0 \\ 0 & (1 - \cos \theta_i) \end{bmatrix} \right), \end{aligned} \quad (30)$$

which concludes the proof. \square

Proof of Lemma 4 As $\sigma^2 \rightarrow 0$, the mean-covariance ratios of both $a_i + b_i$ and $a_i - b_i$ tend to infinity. Therefore, applying Lemma 3 to Eq. (30) (see proof of Lemma 2), we have $(a_i + b_i)(a_i - b_i) \sim \mathcal{N}(\sin^2 \theta_i \alpha_i^2, 4\sigma^2 \sin^2 \theta_i (\alpha_i^2 + \sigma^2))$. Applying the independence between $(a_i + b_i)(a_i - b_i)$ and $(a_j + b_j)(a_j - b_j)$ ($i \neq j$), we obtain the desired result by summing the mean and variance over all i . \square

Proof of Theorem 4 We prove the theorem by deriving upper bounds on $\Pr(\mathcal{C}_2 | \mathcal{C}_1, \boldsymbol{\alpha})$ and $\Pr(\mathcal{C}_1 | \mathcal{C}_2, \boldsymbol{\alpha})$.

$$\begin{aligned} \Pr(\mathcal{C}_2 | \mathcal{C}_1, \boldsymbol{\alpha}) &= \Pr \left(\sum_i (a_i + b_i)(a_i - b_i) \leq 0 \right) \\ &= \Pr \left(\frac{\sum_i (a_i + b_i)(a_i - b_i) - \sum_i \sin^2 \theta_i \alpha_i^2}{2\sigma \sqrt{\sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2)}} \leq \right. \\ &\quad \left. - \frac{\sum_i \sin^2 \theta_i \alpha_i^2}{2\sigma \sqrt{\sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2)}} \right). \end{aligned} \quad (31)$$

As $\sigma \rightarrow 0$, the term to the left of " \leq " in the last line of Eq. (31) is standard normal distributed. Therefore we can invoke the Gaussian tail bound to obtain

$$\begin{aligned} &\Pr(\mathcal{C}_2 | \mathcal{C}_1, \boldsymbol{\alpha}) \\ &= \Pr \left(\frac{\sum_i (a_i + b_i)(a_i - b_i) - \sum_i \sin^2 \theta_i \alpha_i^2}{2\sigma \sqrt{\sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2)}} \geq \right. \\ &\quad \left. - \frac{\sum_i \sin^2 \theta_i \alpha_i^2}{2\sigma \sqrt{\sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2)}} \right) \\ &\leq \frac{1}{2} \exp \left[-\frac{(\sum_i \sin^2 \theta_i \alpha_i^2)^2}{8\sigma^2 \sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2)} \right]. \end{aligned} \quad (32)$$

$\Pr(\mathcal{C}_1 | \mathcal{C}_2, \boldsymbol{\alpha})$ can be upper bounded in the same manner:

$$\Pr(\mathcal{C}_1 | \mathcal{C}_2, \boldsymbol{\alpha}) \leq \frac{1}{2} \exp \left[-\frac{(\sum_i \sin^2 \theta_i \alpha_i^2)^2}{8\sigma^2 \sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2)} \right]. \quad (33)$$

Therefore,

$$\begin{aligned}
P_e &= \frac{1}{2} \int \Pr(\mathcal{C}_2 | \mathcal{C}_1, \boldsymbol{\alpha}) p(\boldsymbol{\alpha}) d\boldsymbol{\alpha} + \frac{1}{2} \int \Pr(\mathcal{C}_1 | \mathcal{C}_2, \boldsymbol{\alpha}) q(\boldsymbol{\alpha}) d\boldsymbol{\alpha} \\
&\leq \int \frac{1}{2} \exp \left[-\frac{(\sum_i \sin^2 \theta_i \alpha_i^2)^2}{8\sigma^2 \sum_i \sin^2 \theta_i (\alpha_i^2 + \sigma^2)} \right] \frac{p(\boldsymbol{\alpha}) + q(\boldsymbol{\alpha})}{2} d\boldsymbol{\alpha} \\
&\triangleq \int \mathcal{E}(\theta, \boldsymbol{\alpha}, \sigma) \frac{p(\boldsymbol{\alpha}) + q(\boldsymbol{\alpha})}{2} d\boldsymbol{\alpha},
\end{aligned} \tag{34}$$

which concludes the proof. \square

APPENDIX E

Proof of Proposition 1 Observe that

$$\|\mathbf{X}^\top \mathbf{P} \mathbf{X} - \mathbf{T}\|_F^2 = \|(\mathbf{X}^\top \otimes \mathbf{X}^\top) \text{vec}(\mathbf{P}) - \text{vec}(\mathbf{T})\|_2^2,$$

is a least squares problem with minimizer

$$\text{vec}(\mathbf{P}^*) = (\mathbf{X}^\top \otimes \mathbf{X}^\top)^\dagger \text{vec}(\mathbf{T}) = \mathbf{X}^\top,$$

which can be rearranged to give

$$\mathbf{P}^* = (\mathbf{X}^\top)^\dagger \mathbf{T} [(\mathbf{X}^\top)^\dagger]^\top = (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X} \mathbf{T} \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \succeq 0.$$

\square

REFERENCES

- [1] R. Basri and D. W. Jacobs, "Lambertian reflectance and linear subspaces," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 2, pp. 218–233, 2003.
- [2] S. R. Rao, R. Tron, R. Vidal, and Y. Ma, "Motion segmentation via robust subspace separation in the presence of outlying, incomplete, or corrupted trajectories," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- [3] J. D. M. Rennie and N. Srebro, "Fast maximum margin matrix factorization for collaborative prediction," in *In Proceedings of the International Conference of Machine Learning*, 2005.
- [4] M. L. McCloud and L. L. Scharf, "A new subspace identification algorithm for high-resolution doa estimation," *IEEE Transactions on Antennas and Propagation*, vol. 50, no. 10, pp. 1382–1390, 2002.
- [5] Y. Xie, J. Huang, and R. Willett, "Change-point detection for high-dimensional time series with missing data," *IEEE Journal of Selected Topics in Signal Processing*, vol. 7, no. 1, pp. 12–27, 2013.
- [6] Y. Chi, Y. C. Eldar, and R. Calderbank, "Petrels: Parallel subspace estimation and tracking using recursive least squares from partial observations," *IEEE Transaction on Signal Processing*, vol. 61, no. 23, pp. 5947–5959, 2013.
- [7] E. Elhamifar and R. Vidal, "Sparse subspace clustering," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [8] Q. Qiu and G. Sapiro, "Learning transformations for clustering and classification," to appear in *Journal of Machine Learning Research*, 2014.
- [9] M. Chen, J. Silva, J. Paisley, C. Wang, D. Dunson, and L. Carin, "Compressive sensing on manifolds using a nonparametric mixture of factor analyzers: Algorithm and performance bounds," *IEEE Transactions on Signal Processing*, vol. 58, no. 12, pp. 6140–6155, 2010.
- [10] S. Ji, Y. Xue, and L. Carin, "Bayesian compressive sensing," *IEEE Transactions on Signal Processing*, vol. 56, no. 6, pp. 2346–2356, 2008.
- [11] G. Yu and G. Sapiro, "Statistical compressed sensing of gaussian mixture models," *IEEE Transactions on Signal Processing*, vol. 59, no. 12, pp. 5842–5858, 2011.
- [12] G. Yu, G. Sapiro, and S. Mallat, "Solving inverse problems with piecewise linear estimators: from gaussian mixture models to structured sparsity," *IEEE Transactions on Image Processing*, vol. 21, no. 5, pp. 2481–2499, 2012.
- [13] F. Renna, R. Calderbank, L. Carin, and M. Rodrigues, "Reconstruction of signals drawn from a gaussian mixture via noisy compressive measurements," *IEEE Transactions on Signal Processing*, vol. 62, no. 9, pp. 2265–2277, 2014.
- [14] M. Chen, W. Carson, M. Rodrigues, R. Calderbank, and L. Carin, "Communications inspired linear discriminant analysis," in *International Conference of Machine Learning*, 2012.
- [15] M. Nokleby, M. Rodrigues, and R. Calderbank, "Discrimination on the grassmann manifold: Fundamental limits of subspace classifiers," *IEEE Transaction on Information Theory*, vol. 61, no. 4, pp. 2133–2147, 2015.
- [16] B. M. Hochwald, T. L. Marzetta, T. J. Richardson, W. Sweldens, and R. Urbanke, "Systematic design of unitary space-time constellations," *IEEE Transactions on Information Theory*, vol. 46, no. 6, pp. 1962–1973, 2000.
- [17] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification (2nd Edition)*. New York, NY: Wiley-Interscience, 2000.
- [18] H. Reboredo, F. Renna, R. Calderbank, and M. Rodrigues, "Compressive classification of a mixture of gaussians: Analysis, designs and geometrical interpretation," 2014, arXiv preprint arXiv:1401.6962.
- [19] A. Edelman, T. A. Arias, and S. T. Smith, "The geometry of algorithms with orthogonality constraints," *SIAM journal on Matrix Analysis and Applications*, vol. 20, no. 2, pp. 303–353, 1998.
- [20] T. Wimalajeewa, H. Chen, and P. K. Varshney, "Performance limits of compressive sensing-based signal classification," *IEEE Transactions on Signal Processing*, vol. 60, no. 6, pp. 2758–2770, 2012.
- [21] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 2, pp. 210–227, 2009.
- [22] A. Bhattacharyya, "On a measure of divergence between two statistical populations defined by their probability distributions," *Bulletin of the Calcutta Mathematical Society*, vol. 45, pp. 99–109, 1943.
- [23] V. Tarokh, N. Seshadri, and R. Calderbank, "Space-time codes for high data rate wireless communication: Performance criterion and code construction," *IEEE Transactions on Information Theory*, vol. 44, no. 2, pp. 744–765, 1998.
- [24] A. Ashikhmin and R. Calderbank, "Grassmannian packings from operator reed-muller codes," *IEEE transactions on Information Theory*, vol. 56, no. 11, pp. 5689–5714, 2010.
- [25] L. A. Aroian, "The probability function of the product of two normally distributed variables," *The Annals of Mathematical Statistics*, pp. 265–271, 1947.
- [26] P. N. Belhumeur, J. P. Hespanha, and D. Kriegman, "Eigenfaces vs. fisherfaces: Recognition using class specific linear projection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 19, no. 7, pp. 711–720, 1997.
- [27] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 2, pp. 210–227, 2009.
- [28] W. Johnson and J. Lindenstrauss, "Extensions of lipschitz mappings into a hilbert space," in *Conference in Modern Analysis and Probability*, 1982.
- [29] T. Sim, S. Baker, and M. Bsat, "The cmu pose, illumination, and expression (pie) database," in *IEEE International Conference on Automatic Face and Gesture Recognition*, 2002.
- [30] D. L. Kleinman and M. Athans, "The design of suboptimal linear time-varying systems," *IEEE Transaction on Automatic Control*, vol. 13, no. 2, pp. 150–159, 1968.
- [31] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge University Press, 2012.