

# Discrimination on the Grassmann Manifold: Fundamental Limits of Subspace Classifiers

Matthew Nokleby, *Member, IEEE*, Miguel Rodrigues, *Member, IEEE*, and Robert Calderbank, *Fellow, IEEE*

## Abstract

We present fundamental limits on the reliable classification of linear and affine subspaces from noisy, linear features. Drawing an analogy between discrimination among subspaces and communication over vector wireless channels, we propose two Shannon-inspired measures to characterize asymptotic classifier performance. First, we define the *classification capacity*, which characterizes necessary and sufficient conditions for the misclassification probability to vanish as the signal dimension, the number of features, and the number of subspaces to be discerned all approach infinity. Second, we define the *diversity-discrimination tradeoff* which, by analogy with the diversity-multiplexing tradeoff of fading vector channels, characterizes relationships between the number of discernible subspaces and the misclassification probability as the noise power approaches zero. We derive upper and lower bounds on these measures which are tight in many regimes. Numerical results, including a face recognition application, validate the results in practice.

## I. INTRODUCTION

The classification of high-dimensional signals arises in a host of situations, from face and digit recognition [1]–[3] to tumor classification [4], [5], and to the music-identification app Shazam [6]. These problems involve massive data sets—images with millions of pixels, DNA arrays with thousands of genes, or audio clips with tens of thousands of samples—which presents a substantial burden of computation and storage. Frequently, however, the data lie near a low-dimensional subspace of ambient space. For example, images of an individual’s face, subject to constraints on pose, lighting, and convexity, lie almost entirely on a subspace of five to nine dimensions, regardless of the ambient dimension of the image [7]–[9]. One therefore can pose classification tasks like face recognition as subspace classification problems.

When identifying low-dimensional subspaces, one can reduce the computation and storage burden by classifying from a low-dimensional representation of the signal of interest. This process is called *feature extraction*, and standard techniques, including linear discriminant analysis (LDA) and principal component analysis (PCA) [10], as well as their myriad variations, are well studied. One pays a price, however, for computational tractability. In principle, extracting low-dimensional features from high-dimensional signals degrades classifier performance, and it is unclear *a priori* how many features are necessary to ensure success.

In this paper, we present a rigorous, information-theoretic characterization of classifier performance of high-dimensional data from low-dimensional features. We show that performance depends on several factors, including the number of subspaces to be discriminated, the number of features extracted, and the underlying subspace structure. In particular, we consider the classification of  $k$ -dimensional linear and affine subspaces of  $\mathbb{R}^N$  from  $M$  linear features corrupted by Gaussian noise. To characterize classifier performance, we define two performance measures:

- The *classification capacity*, which characterizes the number of unique subspaces that can be discerned as a function of the noise power,  $N$ ,  $M$ , and  $k$ , as the latter three quantities approach infinity. Just as the usual Shannon capacity captures the phase transition of the error probability, as a function of the information rate, as the code length goes to infinity, the classification capacity captures the phase transition of the misclassification probability, in terms of the (logarithm of the) number of subspaces, as the signal dimension goes to infinity.
- The *diversity-discrimination tradeoff* (DDT), which characterizes the relationship between the number of subspaces and the misclassification probability as the noise power goes to zero. Just as the diversity-multiplexing

This work supported in part by AFOSR grant FA9550-13-1-0076 and the Royal Society International Exchanges Scheme IE120996. Preliminary elements of this work were presented at the IEEE Information Theory Workshop, Sept. 2013 and the IEEE Symposium on Information Theory, July 2014. Matthew Nokleby and Robert Calderbank are with Duke University, Durham, NC (emails: {matthew.nokleby,robert.calderbank}@duke.edu), and Miguel Rodrigues is with University College London (email: m.rodrigues@ucl.ac.uk)

tradeoff for fading vector wireless channels [11] characterizes the number of codewords and the error probability in terms of a region of achievable exponent pairs in the signal-to-noise ratio (SNR), the DDT specifies a region of achievable exponent pairs in the noise power for the number of subspaces and the misclassification probability.

The motivation for the preceding definitions is an analogy between classification from noisy features and communication over non-coherent vector channels. Indeed, the title of our paper alludes to [12], which investigates the capacity of the block-fading non-coherent channel in geometric terms. It shows that, at high SNR and for sufficiently long coherence time, transmitters achieve near-capacity rates by sending *subspaces* as codewords. Therefore the decoding task is to discern subspaces from noisy observations, and the capacity corresponds to asymptotic packings in the Grassmann manifold. Further works give tighter bounds on the capacity and explore the diversity-multiplexing tradeoff of the non-coherent channel [12]–[15].

For the classification of  $k$ -dimensional linear subspaces from noisy, linear features, one can demonstrate a syntactic duality with communications over non-coherent vector channels. Specifically, the classification problem is dual to a non-coherent communications over a channel with  $k$  transmit antennas, a single receive antenna, and a coherence time of  $M$ . In a preliminary version of this work, we applied results from [12], [14] directly to prove necessary conditions for successful classification [16]. These bounds translate into upper bounds on the classification capacity and diversity-discrimination tradeoff considered in this paper.

However, these bounds are somewhat crude. In the dual communications problem, the optimal transmission strategy employs only a single transmit antenna, which is equivalent to classifying subspaces of dimension  $k = 1$ . Therefore, the upper bounds are loose when classifying subspaces of higher dimension. Furthermore, because the classification problem is not known to be *information stable* [17], the mutual information between subspaces and features does not lower bound the classification capacity even for  $k = 1$ . To prove tighter upper bounds on performance, we develop new bounds on the mutual information, and to prove lower bounds on performance we analyze the misclassification probability directly.

### A. Summary of Results

Our primary contributions are upper and lower bounds, which are tight in many regimes, on the classification capacity and the diversity-discrimination tradeoff.

In Section III, we study the classification capacity. First we consider the classification of linear subspaces, which we model by taking the classes to follow zero-mean Gaussian distributions with approximately low-rank covariances. The covariances have two components: a rank  $k$  component corresponding to the class subspace, and an identity component scaled by  $\sigma^2$  corresponding to deviations from the subspace. We further suppose a prior distribution the subspaces which is uniform over the Grassmann manifold of  $k$ -dimensional subspaces in  $\mathbb{R}^N$ . We present an upper bound on the classification capacity, showing almost surely that the number of subspaces cannot scale any faster than  $(1/\sigma^2)^{\frac{M-k}{2}}$ . This result is intuitive: The lower the inherent signal dimension, the fewer features are required to classify the signal reliably. We also present a lower bound on the classification capacity, showing that the misclassification probability decays to zero, except for a set of subspaces having vanishing probability, provided the number of subspaces grows slower than  $(1/\sigma^2)^{\frac{\min\{k, M-k\}}{2}}$ . When  $M \leq 2k$ , the bounds are tight up to a  $O(1)$  term. Furthermore, based on simulations presented in Section V, we conjecture that the upper bound is tight and that the gap between lower and upper bounds when  $M > 2k$  is merely an artifact of the analysis.

We then consider the classification of *affine* subspaces, or linear subspaces translated by nonzero points. We model affine spaces by taking the classes to again be modeled by approximately rank- $k$  covariances, but this time to have nonzero means. We again suppose a uniform prior over the Grassmann manifold, and we further suppose that the means are distributed according to a standard Gaussian distribution. We characterize the classification capacity up to a  $O(1)$  term, showing that the number of subspaces growing no faster than  $(1/\sigma^2)^{\frac{M-k}{2}}$  is both necessary and sufficient for the probability of error to decay to zero, again except for a set of subspaces having vanishing probability.

In Section IV, we study the diversity-discrimination tradeoff. For linear subspaces, we derive an upper bound, showing that the average misclassification probability decays no faster than  $(1/\sigma^2)^{-\frac{\min\{k, M-k\}}{2}}$  as  $\sigma^2 \rightarrow 0$  and that the misclassification probability exhibits an error floor when the number of subspaces scales faster than  $(1/\sigma^2)^{\frac{M-k}{2}}$ . We also derive a lower bound, showing that the misclassification capacity decays at least as  $(1/\sigma^2)^{-\frac{\min\{k, M-k\}-r}{2}}$ .

when the number of subspaces scales as  $(1/\sigma^2)^{\frac{r}{2}}$ . For affine spaces, we specify the DDT exactly, showing that the average probability decays as  $(1/\sigma^2)^{-\frac{M-k-r}{2}}$  when the number of subspaces scales as  $(1/\sigma^2)^{\frac{r}{2}}$ .

For both linear and affine subspace classification, the lower bounds on performance are realized by any feature matrix having  $M$  orthonormal rows in  $\mathbb{R}^N$ . Therefore, for regimes in which the bounds are tight, the asymptotic performance as characterized by the classification capacity and DDT is invariant to rotations of the linear features.

In Section V, we evaluate our claims empirically. We first examine the error performance of classifiers over randomly-drawn linear subspaces, focusing on the regimes in which the upper and lower bounds disagree. Then, we test the correspondence of our theoretical results to a practical face recognition application. Using standard classification algorithms against public datasets, we observe error performance that agrees with our predictions to within a reasonable tolerance.

## B. Prior Work

The statistics and machine learning literature contains a large body of work on feature extraction or supervised dimensionality reduction. In addition to the venerable linear discriminant analysis and principal component analysis, which depend only on the second-order statistics of the data, linear techniques based on higher-order statistics were developed in [18]–[26]. Owing to Fano’s inequality, the algorithm of [18] chooses linear features having maximal the mutual information with the classes, whereas [20], [25], [26] employ approximations to the mutual information based on Rényi entropy. In [27] linear features are chosen for subspace classification according to a nuclear-norm optimization problem, and in [28] an LDA-inspired Grassmann discriminant analysis is proposed. Finally, nonlinear dimensionality reduction techniques have recently become popular [29], [30].

In the signal processing literature, information-theoretic limits on subspace classification arise under the framework of sparse support recovery. The set of all  $k$ -sparse vectors in  $\mathbb{R}^N$  is a union of subspaces, and recovering the sparsity pattern is equivalent to finding the subspace in which the signal lies. A (data) deluge of recent works [31]–[43] provides necessary and sufficient scaling laws on the triplet  $(N, k, M)$ , where  $M$  is the number of compressive measurements taken, for recovery of a sparse signal. Different assumptions on the measurement matrices, decoders, error metrics, and sparsity regimes give rise to different scaling laws. While these works do provide fundamental limits on subspace classifier performance, sparse support recovery entails a specialization to the union of canonical subspaces, and the results presented in the preceding works do not bear directly on our study.

Reference [44] considers compressed learning, i.e. learning directly in the compressive measurement domain rather than in the original data domain, showing that when data admit a sparse representation, low-dimensional feature extraction preserves the learnability and the separability of the data. Along a similar vein, a recent work [45] considers the compressive classification of convex sets, proving limits on the number of measurements required to ensure that the projected sets remain separated.

A few works have focused on the classification of Gaussian mixtures, which is closely related to the linear and affine subspace classification considered herein. In [46] classifier performance is studied for a finite number of classes as a function of signal geometry; these results prefigure the DDT results presented in the sequel. In [47], the number of measurements required to classify and reconstruct a signal drawn from a Gaussian mixture is characterized.

Researchers have also studied information-theoretic limits on other classification problems. The authors of [48] provide asymptotic limits on the success of model selection of Markov random fields. The authors of [49] use results in universal source coding to prove general bounds on classifier performance. The authors of [50] study the limits of database recovery from low-dimensional features, characterizing an “identification capacity” which is analogous to the classification capacity studied in this paper.

## C. Notation

We let bold lowercase letters denote vectors and bold uppercase letters to denote matrices. We let  $\mathbb{R}$  and  $\mathbb{Z}$  denote the field of reals and integers, respectively. We let  $\mathbf{I}$  and  $\mathbf{0}$  denote the identity matrix and the all-zeros matrix, respectively, indicating the size of the matrix in a subscript when necessary. We let  $\|\cdot\|$  denote the Euclidean norm; when applied to a matrix it denotes the induced operator norm. We let  $E[\cdot]$  denote the expectation, indicating the distribution over which the expectation is taken by a subscript when necessary. We let  $[\cdot]$  and  $[\cdot]^+$  denote the floor function and the positive part of a number, respectively. We let  $\text{eig}(\cdot)$  denote the vector of eigenvalues of a square

matrix. We let  $\stackrel{d}{=}$  denote equality in distribution. We let  $\mathcal{N}(\mu, \Sigma)$  denote a Gaussian distribution with mean  $\mu$  and covariance matrix  $\Sigma$ . We let  $\mathcal{W}_M(N, \mathbf{V})$  denote the  $M \times M$  Wishart distribution with degrees of freedom  $N$  and shape matrix  $\mathbf{V}$ .

## II. PRELIMINARIES

### A. Problem Definition

We consider the statistical classification problem, in which the signal of interest  $\mathbf{x} \in \mathbb{R}^N$  is distributed according to one of  $L$  class-conditional densities  $p_l(\mathbf{x})$ , each of which is known to the classifier. The classifier observes noisy linear projections of  $\mathbf{x}$ , from which it attempts to determine the class-conditional density from which  $\mathbf{x}$  was drawn. These projections, denoted by  $\mathbf{y} \in \mathbb{R}^M$ , are related to the signal  $\mathbf{x} \in \mathbb{R}^N$  as follows:

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{z}, \quad (1)$$

where  $\Phi \in \mathbb{R}^{M \times N}$  is a matrix describing the linear features, and  $\mathbf{z} \in \mathbb{R}^M$  is white Gaussian noise with mean zero and per-component variance  $\sigma^2$  for some  $\sigma^2 > 0$ . We suppose  $M \leq N$ , and we constrain  $\|\Phi\| \leq 1$ . The noise  $\mathbf{z}$  describes the deviation between the postulated subspace model and the true signal of interest.<sup>1</sup> Signals will not lie perfectly on the linear or affine subspaces, so we suppose that the projected signal lies approximately within a ball centered at the specified subspaces and having radius  $\sqrt{M}\sigma$ .

To model the classification of linear and affine subspaces, we impose structure on the class-conditional densities  $p_l(\mathbf{x})$ . In particular, we suppose that the class conditional densities are Gaussian with low-rank covariances that correspond to the subspaces. In the case of linear subspaces, these Gaussians have zero mean. In the case of affine subspaces, which are simply translations of linear subspaces, the Gaussians have nonzero means.

We therefore define two sets. For the classification of  $k$ -dimensional linear subspaces of  $\mathbb{R}^N$ , the class-conditional densities  $p_l(\mathbf{x})$  belong to<sup>2</sup>

$$\mathcal{Q}_{\text{linear}}(N, k) = \{\mathcal{N}(\mathbf{0}, \mathbf{U}\mathbf{U}^T) : \mathbf{U} \in \mathbb{R}^{N \times k}\}. \quad (2)$$

In other words, each class-conditional density is a Gaussian supported on the  $k$ -dimensional subspace spanned by the columns of  $\mathbf{U}$ . Similarly, for the classification of  $k$ -dimensional affine spaces, the class-conditional densities belong to

$$\mathcal{Q}_{\text{affine}}(N, k) = \{\mathcal{N}(\mu, \mathbf{U}\mathbf{U}^T) : \mu \in \mathbb{R}^N, \mathbf{U} \in \mathbb{R}^{N \times k}\}. \quad (3)$$

That is, the class-conditional densities are supported on a  $k$ -dimensional subspace as before, but here they are translated by a non-zero vector  $\mu$ .

We parameterize the sets  $\mathcal{Q}_{\text{linear}}(N, k)$  and  $\mathcal{Q}_{\text{affine}}(N, k)$  by the following two sets

$$\mathcal{A}_{\text{linear}}(N, k) = \mathbb{R}^{N \times k} \quad (4)$$

$$\mathcal{A}_{\text{affine}}(N, k) = \mathbb{R}^N \times \mathbb{R}^{N \times k}. \quad (5)$$

Clearly  $\mathcal{A}_{\text{linear}}(N, k)$  and  $\mathcal{A}_{\text{affine}}(N, k)$  are isomorphic to  $\mathcal{Q}_{\text{linear}}(N, k)$  and  $\mathcal{Q}_{\text{affine}}(N, k)$ , respectively. We can represent a linear or affine subspace classification problem by a tuple  $\mathbf{a} = (a_1, \dots, a_L) \in \mathcal{A}^L(N, k)$ , where  $\mathcal{A}^L$  is the  $L$ -fold Cartesian product of  $\mathcal{A}(N, k)$ . The tuple  $\mathbf{a}$  encodes the  $L$  covariances and, when appropriate, the  $L$  means corresponding to the subspaces to be classified. Let  $p(\mathbf{x}; a_l) = p_l(\mathbf{x})$ , for  $1 \leq l \leq L$ , denote the class-conditional densities parameterized by  $\mathbf{a} \in \mathcal{A}(N, k)$ .

Let  $\hat{l} = f(\mathbf{y})$  denote the classifier output, where  $f$  is a mapping from  $\mathbb{R}^M$  to  $\{1, \dots, L\}$ . Then, for a classification problem described by the tuple  $\mathbf{a}$ , define the average misclassification probability:

$$P_\epsilon(\mathbf{a}) = \min_{\|\Phi\| \leq 1} \frac{1}{L} \sum_{l=1}^L \Pr(\hat{l} \neq l | \mathbf{x} \sim p(\mathbf{x}; a_l)), \quad (6)$$

<sup>1</sup>Equivalently, we could remove the additive noise  $\mathbf{z}$  and add a  $\sigma^2 \mathbf{I}$  term to the covariance matrix of each class-conditional density.

<sup>2</sup>We will drop the subscripts linear and affine throughout when discussing classification generally rather than particularizing to linear or affine subspaces.

where each term in the sum is the misclassification probability when  $\mathbf{x}$  is drawn according to  $p(\mathbf{x}; a_l)$ . Observe that we define  $P_e(\mathbf{a})$  in terms of the best feature matrix  $\Phi$ . Therefore, in proving our results we will characterize the feature matrix that achieves optimal classifier performance in the asymptote.

The focus of this paper is the analysis of  $P_e$  in two asymptotic regimes: (i) as the signal dimensions  $N, M, k$  go to infinity, and (ii) as the noise power  $\sigma^2$  goes to zero. In the first case, we derive conditions for which the probability of error decays to zero, except for a set of vanishing probability over  $\mathcal{A}^L$ . In the second case, we derive scaling laws on the probability of error, *averaged* over the possible choices of  $\mathbf{a} \in \mathcal{A}^L$ . To this end, define the following probability distributions over the parameter sets  $\mathcal{A}_{\text{linear}}(N, k)$  and  $\mathcal{A}_{\text{affine}}(N, k)$ :

$$p_{\text{linear}}(a) = \prod_{i=1}^N \prod_{j=1}^k \mathcal{N}(u_{ij}; 0, 1/k) \quad (7)$$

$$p_{\text{affine}}(a) = \prod_{i=1}^N \prod_{j=1}^k \mathcal{N}(u_{ij}; 0, 1/k) \cdot \prod_{l=1}^N \mathcal{N}(\mu_l; 0, 1), \quad (8)$$

where  $u_{ij}$  is the  $(i, j)$ th element of the matrix  $\mathbf{U}$  and  $\mu_i$  is the  $i$ th element of the vector  $\mu$ . These distributions define a measure over the sets of class-conditional densities. In other words, in computing probabilities we suppose that the elements of the matrix  $\mathbf{U}$  and the mean vector  $\mu$  are standard i.i.d. Gaussian.

For both  $p_{\text{linear}}$  and  $p_{\text{affine}}$ , the distribution is supported over the entire parameter space, is invariant to rotations, and yields finite expected signal energy. Specifically, the distribution over the bases  $\mathbf{U}$  is isotropic, which means that the linear subspaces are drawn uniformly from the Grassmann manifold. Therefore, our analysis characterizes classifier performance when “nature” presents us with subspaces without favoring a particular region of the Grassmann; we contend that this assumption is reasonable. Furthermore, while changes to the distributions  $p_{\text{linear}}$  and  $p_{\text{affine}}$  will change the classification capacity and DDT in general, the coarse behavior is robust to variations. In particular, one can recover our proofs subject to straightforward constraints on the eigenvalue distribution of  $\mathbf{U}^T \mathbf{U}$ , showing bounds on the classification capacity that differ at most by a  $O(1)$  term and DDT bounds that agree exactly.

Next, we define the classification capacity and the diversity-discrimination tradeoff.

### B. Classification Capacity

The classification capacity characterizes fundamental performance limits as the signal dimensions approach infinity. We derive bounds on how fast the number of subspaces can grow, as a function of  $N$ ,  $M$ , and  $k$ , while ensuring the misclassification probability decays to zero almost surely.

By analogy with the sequence of codebooks defined for the Shannon capacity, we characterize the classification capacity in terms of a sequence of classification problems indexed by  $M$ . We let the number of features  $M$  grow to infinity, and we let the dimensions  $N$  and  $k$  scale linearly with  $M$  as follows:

$$N(M) = \lfloor \nu M \rfloor, k(M) = \lfloor \kappa M \rfloor, \quad (9)$$

for  $\nu \geq 1$  and  $0 < \kappa < 1$ . We also let the number of subspaces  $L$  scale exponentially in  $M$  as follows:

$$L(M) = \lfloor 2^{\rho M} \rfloor, \quad (10)$$

for some  $\rho \geq 0$ . By analogy with communications theory, the quantity  $\rho$  can be interpreted as the “rate” of the sequence of class alphabets, or the average number of bits discerned per feature if classification succeeds. Indeed, in the sequel we refer to  $\rho$  as the *classification rate*.

*Definition 1:* Fix the dimension ratios  $\nu$  and  $\kappa$  and the classification rate  $\rho$ . Then, define the set of classification problems for which the probability of classification error exceeds an arbitrary small constant  $\epsilon > 0$ :

$$\mathcal{E}(M) = \{\mathbf{a} \in \mathcal{A}^{L(M)}(N(M), k(M)) : P_e(\mathbf{a}) > \epsilon\}. \quad (11)$$

Then, we say that  $\rho$  is *achievable* provided

$$\lim_{M \rightarrow \infty} \int_{\mathcal{E}(M)} \prod_{i=1}^{L(M)} p(a_i) d\mathbf{a} = 0, \quad (12)$$

for any fixed  $\epsilon > 0$ .

Observe that a classification rate  $\rho$  is achievable if

$$\lim_{M \rightarrow \infty} E[\mathbb{P}_\epsilon(\mathbf{a})] = \lim_{M \rightarrow \infty} \int_{\mathcal{A}^{L(M)}(N(M), k(M))} \mathbb{P}_\epsilon(\mathbf{a}) \prod_{i=1}^{L(M)} p(a_i) d\mathbf{a} = 0. \quad (13)$$

This observation follows by contradiction: If there is a subset of  $\mathcal{A}^L(N, k)$  having non-trivial probability for which the misclassification probability remains bounded away from zero, the expected error also remains bounded away from zero.

*Definition 2:* Fix the dimension ratios  $\nu$  and  $\kappa$ . The *classification capacity*, denoted by  $C_{\text{linear}}(\nu, \kappa)$  and  $C_{\text{affine}}(\nu, \kappa)$  for linear and affine space classification, respectively, is the supremum over achievable classification rates  $\rho$ .

In other words, if the classification rate is smaller than  $C(\nu, \kappa)$ , then the probability of classification error approaches zero almost surely over the set of subspace classification problems. Otherwise, the error probability remains bounded away from zero for a non-trivial subset of  $\mathcal{A}^L(N, k)$ .

Although the classification capacity is defined to characterize classifier behavior when  $N$  and  $k$  scale linearly in  $M$  and  $L$  scales exponentially in  $M$ , it also captures other regimes automatically. For example, if  $k$  scales sub-linearly in  $M$ , the asymptotic behavior is the same as if  $\kappa \rightarrow 0$ . Similarly, if  $L$  scales sub-exponentially in  $M$ , the misclassification probability decays to zero whenever the classification capacity is nonzero, and if  $L$  scales super-exponentially the misclassification capacity remains bounded from zero whenever the classification capacity is finite. In view of Theorems 1 and 2, this implies that whenever  $\kappa > 0$  and the number of subspaces grows polynomially in  $M$ , the misclassification probability goes to zero. Because we are dealing with subspaces, it is impossible to have  $N$  scale sub-linearly or  $k$  scale super-linearly in  $M$ . However, at least one regime remains unspecified by our analysis: If  $N$  scales super-linearly in  $M$ , classifier behavior is unclear.

We can bound the classification capacity via the mutual information between the vector  $\mathbf{a} \in \mathcal{A}$  and the feature vector  $\mathbf{y}$ .

*Lemma 1:* The classification capacity satisfies

$$C \leq \lim_{M \rightarrow \infty} \max_{\|\Phi\| \leq 1} \frac{I(\mathbf{a}; \mathbf{y})}{M}, \quad (14)$$

where the mutual information is computed with respect to  $p_{\text{linear}}(a)$  or  $p_{\text{affine}}(a)$  as appropriate.

*Proof:* The proof follows from Fano's inequality. By the standard arguments (e.g. from [51]), we obtain

$$P_e(\mathbf{a}) \geq 1 - \frac{\max_{\|\Phi\| \leq 1} I(\mathbf{a}; \mathbf{y}) - 1}{M\rho},$$

which is bounded away from zero when  $\rho$  exceeds the RHS of (14). ■

Observe that we have proven more than just an upper bound on the capacity. When  $\rho$  exceeds RHS of (14), not only is there a non-trivial set for which the error probability remains positive, but that set is also all of  $\mathcal{A}^L(N, k)$ . If the upper bound of Lemma 1 is tight, then the mutual information characterizes a sharp phase transition in the error probability. If the number of subspaces grows sufficiently slowly, the probability of error vanishes almost everywhere; otherwise, is bounded away from zero everywhere.

However, it is not clear whether Lemma 1 is tight. If the ‘‘channel’’ between subspaces and features is *information stable*—meaning roughly that the normalized information density converges on the normalized mutual information—then the mutual information completely characterizes the classification capacity and (14) holds with equality [17]. Alternatively, applying the results of [52], one can express the classification capacity directly in terms of the information density. Analysis of the information density is difficult, however, as is the verification of information stability, so Lemma 1 remains an upper bound only.

To prove lower bounds on the classification capacity, we analyze directly the misclassification probability. Our main tool is the Bhattacharyya bound on the pairwise misclassification probability [53], [54], which we state here for Gaussian distributions.

*Lemma 2:* Suppose we observe a signal that is distributed according to  $\mathcal{N}(\mu_1, \Sigma_1)$  or  $\mathcal{N}(\mu_2, \Sigma_2)$  with equal prior

probability. Define

$$B = \frac{1}{2} \ln \left( \frac{|\frac{\Sigma_1 + \Sigma_2}{2}|}{|\Sigma_1|^{\frac{1}{2}} |\Sigma_2|^{\frac{1}{2}}} \right) + \frac{1}{8} (\mu_1 - \mu_2) \left[ \frac{\Sigma_1 + \Sigma_2}{2} \right]^{-1} (\mu_1 - \mu_2). \quad (15)$$

Then, supposing maximum likelihood classification, the misclassification probability is bounded by

$$P_e((\mu_1, \Sigma_1, \mu_2, \Sigma_2)) \leq \frac{1}{2} \exp(-B). \quad (16)$$

In [54] it is also observed that the Bhattacharyya bound is exponentially tight in the sense that, if the pairwise error decays to zero, it approaches  $c \cdot \exp(-B)$  for some constant  $c$ . A consequence of this observation, which we will see in Section IV, is that the Bhattacharyya bound predicts the maximum diversity gain for both linear and affine subspace classifiers.

### C. Diversity-Discrimination Tradeoff

The *diversity-multiplexing tradeoff* (DMT) was introduced in the context of wireless communications to characterize the high-SNR performance of fading vector channels. It was shown in [11] that the spatial flexibility provided by multiple antennas can simultaneously increase the achievable rate and decrease the probability of error, but only according to a tradeoff that is precisely characterized at high SNR. We define a similar characterization in the context of classification, called the *diversity-discrimination tradeoff* (DDT), which captures the relationship between the increase of discernible subspaces and the decay of misclassification probability as the noise power approaches zero.

For the DDT, we keep  $N$ ,  $M$ , and  $k$  fixed, but we let the number of subspaces scale in the noise power as follows:

$$L(\sigma^2) = \lfloor (1/\sigma^2)^{\frac{r}{2}} \rfloor, \quad (17)$$

for some  $r \geq 0$ , which we call the *discrimination gain*. We define the DDT in terms of the misclassification probability averaged over the ensemble of classification problems, which we denote by

$$\bar{P}_e(\sigma^2, r) = E[P_e(\mathbf{a})] = \int_{\mathcal{A}(N,k)} P_e(\mathbf{a}) \prod_{l=1}^{L(\sigma^2)} p(a_l) d\mathbf{a}, \quad (18)$$

where here we express the probability of error as a function of the discrimination gain  $r$  and the noise power  $\sigma^2$ . Specifically, the diversity-discrimination tradeoff is defined as the following function:

$$d(r) = \lim_{\sigma^2 \rightarrow 0} -\frac{\log \bar{P}_e(\sigma^2, r)}{\frac{1}{2} \log(1/\sigma^2)}. \quad (19)$$

We refer to  $d(r)$  as the *diversity gain* for discrimination gain  $r$ . In other words, when the number of subspaces increases as  $(1/\sigma^2)^{r/2}$ , the probability of error decays as  $(1/\sigma^2)^{-d(r)/2} + o(\log(\sigma^2))$ . In the sequel we refer to  $d_{\text{linear}}(r)$  and  $d_{\text{affine}}(r)$  as appropriate.

By contrast to the classification capacity, where we characterize phase transitions in the error probability that hold almost surely, for the DDT we specify scaling laws in the error probability that hold on the average over  $\mathcal{A}$ . Rather than specifying *if* the probability of error decays to zero, the DDT specifies *how quickly* it decays. In the former case, it is straightforward to define the failure event and show that it has vanishing probability. In the latter case, it is unclear how to define such a failure event, so we state only an average-case result.

As with the classification capacity, we can derive bounds on the DDT from the mutual information.

*Lemma 3:* Fix  $N$ ,  $M$ , and  $k$ . Then,  $d(r) = 0$  whenever

$$r \geq \lim_{\sigma^2 \rightarrow 0} \max_{\Phi, \|\Phi\| \leq 1} \frac{I(\mathbf{a}; \mathbf{y})}{\frac{1}{2} \log(1/\sigma^2)}, \quad (20)$$

where again the mutual information is calculated with respect to  $p_{\text{linear}}(a)$  or  $p_{\text{affine}}(a)$  as appropriate.

*Proof:* Again we invoke Fano's inequality. Whenever  $r$  is as large as the specified quantity, the probability of error is bounded away from zero, and the diversity gain is zero by definition. ■

### III. CLASSIFICATION CAPACITY

Here we characterize the classification capacities  $C_{\text{linear}}(\nu, \kappa)$  and  $C_{\text{affine}}(\nu, \kappa)$ . We prove upper bounds that show that, for both linear and affine subspace classification, the probability of error remains bounded away from zero almost surely whenever the number of subspaces scales faster than  $(1/\sigma^2)^{\frac{M-k}{2}}$ . For linear spaces, we prove a lower bound which matches the upper bound to within an  $O(1)$  term for  $\kappa \geq 1/2$ ; otherwise the bounds disagree. For affine spaces, we prove a lower bound which is tight to within an  $O(1)$  term for all  $\kappa$ . This suggests the somewhat surprising conclusion that, at least for  $\kappa \geq 1/2$ , translating subspaces by nonzero vectors does not substantially increase the number of subspaces a classifier can discriminate. Whether this conclusion extends to  $\kappa < 1/2$  depends on the tightness of the upper bound. However, as we will see in Section IV, affine subspaces are easier to discriminate in the sense that the misclassification probability decays faster as  $\sigma^2 \rightarrow 0$ .

#### A. Linear Subspaces

First, we bound on  $C_{\text{linear}}(\nu, \kappa)$ .

*Theorem 1:* For linear subspace classification, the classification capacity is bounded by

$$\frac{\min\{\kappa, 1 - \kappa\}}{2} \log_2 \left( 1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2}{\sigma^2} \right) - \frac{\kappa}{2} \leq C_{\text{linear}}(\nu, \kappa) \leq \frac{1 - \kappa}{2} \log_2 \left( \frac{1}{\sigma^2} \right) + \frac{1}{2} \log_2(1 + \sigma^2) - \frac{\kappa}{2} \log_2((\sqrt{1/\kappa} - 1)^2 + \sigma^2). \quad (21)$$

*Proof:* We first prove the upper bound by estimating the mutual information between the subspaces and the features and invoking Lemma 1. Then, we prove the lower bound by invoking Lemma 2 and applying the union bound.

**Upper Bound:** To bound the mutual information  $I(\mathbf{a}; \mathbf{y}) = I(\mathbf{U}; \mathbf{y})$ , we first characterize the optimum choice of  $\Phi$ . Following the argument in [18, Theorem 2], we compute the gradient of the mutual information with respect to the singular values of  $\Phi$ . Writing the singular value decomposition as  $\Phi = \mathbf{W}_\Phi \Lambda_\Phi \mathbf{V}_\Phi^T$ , the gradient is

$$\nabla_{\Lambda_\Phi} I(\mathbf{U}; \mathbf{y}) = \Lambda_\Phi \mathbf{V}_\Phi^T E \left[ \int p(\mathbf{y}|\mathbf{U})(\mathbf{m} - \mathbf{m}_\mathbf{U})(\mathbf{m} - \mathbf{m}_\mathbf{U})^T d\mathbf{y} \right] \mathbf{V}_\Phi, \quad (22)$$

where

$$\mathbf{m} = \int \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

is the mean with respect to the posterior distribution, and

$$\mathbf{m}_\mathbf{U} = \int \mathbf{x} p(\mathbf{x}|\mathbf{y}, \mathbf{U}) d\mathbf{x}$$

is the mean with respect to the conditional posterior. Observe from (22) that the diagonal elements of the gradient are non-negative, which implies that the mutual information is non-decreasing with the singular values of  $\Phi$ . Because we constrain  $\|\Phi\| \leq 1$ , it follows that the singular values of the optimal  $\Phi$  are identically unity.

Assuming this condition on  $\Phi$ , we bound the mutual information. By definition,

$$I(\mathbf{U}; \mathbf{y}) = h(\mathbf{y}) - h(\mathbf{y}|\mathbf{U}).$$

To bound the conditional entropy, observe that the conditional distribution of  $\mathbf{y}$  is

$$p(\mathbf{y}|\mathbf{U}) = \mathcal{N}(0, \Phi \mathbf{U} \mathbf{U}^T \Phi^T + \sigma^2 \cdot \mathbf{I}).$$

Let  $\lambda_i$  denote the  $i$ th ordered eigenvalue of  $\mathbf{U}^T \Phi^T \Phi \mathbf{U}$ . Then, the conditional entropy is

$$\begin{aligned} h(\mathbf{y}|\mathbf{U}) &= \sum_{i=1}^k \frac{1}{2} E[\log_2(2\pi e(\lambda_i + \sigma^2))] + \frac{M-k}{2} \log_2(2\pi e\sigma^2) \\ &\geq \frac{k}{2} E[\log_2(\lambda_k + \sigma^2)] + \frac{M-k}{2} \log_2(\sigma^2) + \frac{M}{2} \log_2(2\pi e), \end{aligned}$$



where the expectation is with respect to  $\mathbf{U}$ , and where the inequality is trivially obtained by substituting the smallest positive eigenvalue  $\lambda_k$ . We next bound this eigenvalue. Because  $\Phi$  has singular values identically equal to unity, and because  $\mathbf{U}$  has i.i.d. Gaussian entries with variance  $1/k$ , the matrix  $\Phi\mathbf{U} \in \mathbb{R}^{M \times k}$  also has i.i.d. Gaussian entries with variance  $1/k$ . Therefore,

$$\mathbf{U}^T \Phi^T \Phi \mathbf{U} \stackrel{d}{=} \frac{1}{k} \mathbf{W},$$

where  $\mathbf{W} \sim \mathcal{W}_k(\mathbf{I}, M)$ . In [55, Theorem 1] it is shown that the smallest eigenvalue of  $1/M \cdot \mathbf{W}$  converges to  $(1 - \sqrt{\kappa})^2$  almost surely as  $M \rightarrow \infty$ . Therefore, the minimum eigenvalue of  $\mathbf{U}^T \Phi^T \Phi \mathbf{U}$ , which is equal to  $\lambda_k$ , converges on  $(\sqrt{1/\kappa} - 1)^2$  almost surely. We can therefore bound the conditional mutual information by

$$h(\mathbf{y}|\mathbf{U}) \geq \frac{k}{2} \log_2((\sqrt{1/\kappa} - 1)^2 + \epsilon(M) + \sigma^2) + \frac{M-k}{2} \log_2(\sigma^2) + \frac{M}{2} \log_2(2\pi e), \quad (23)$$

where  $\epsilon(M) \rightarrow 0$  as  $M \rightarrow \infty$ .

We turn next to the differential entropy of  $\mathbf{y}$ . We first compute the expected covariance, which is

$$\begin{aligned} E[\mathbf{y}\mathbf{y}^T] &= E_{\mathbf{U}}[\Phi\mathbf{U}\mathbf{U}^T\Phi + \sigma^2\mathbf{I}] \\ &= (1 + \sigma^2)\mathbf{I}. \end{aligned}$$

Noting that the differential entropy for a fixed covariance is maximized by the multivariate Gaussian distribution, we obtain

$$h(\mathbf{y}) \leq \frac{M}{2} \log_2(2\pi e(1 + \sigma^2)). \quad (24)$$

Combining terms and letting  $M \rightarrow \infty$ , we finally obtain

$$\lim_{M \rightarrow \infty} \max_{\Phi} \frac{I(\mathbf{y}; \mathbf{U})}{M} \leq \frac{1-\kappa}{2} \log\left(\frac{1}{\sigma^2}\right) + \frac{1}{2} \log_2(1 + \sigma^2) - \frac{\kappa}{2} \log_2((\sqrt{1/\kappa} - 1)^2 + \sigma^2). \quad (25)$$

Applying Lemma 1 to (25), we obtain the upper bound.

**Lower Bound:** Choose  $\Phi \in \mathbb{R}^{M \times N}$  to be any matrix with orthonormal rows. Observe that while this choice maximizes the mutual information, it does not minimize the probability of error in general. Applying the Bhat-tacharyya bound from Lemma 2, the probability of a pairwise error between two subspaces  $i$  and  $j$  is bounded by

$$P_e(\mathbf{U}_i, \mathbf{U}_j) \leq \frac{1}{2} \cdot \left( \frac{|\frac{\Phi\mathbf{U}_i\mathbf{U}_i^T\Phi^T + \Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T + 2\sigma^2\mathbf{I}}{2}|}{|\Phi\mathbf{U}_i\mathbf{U}_i^T\Phi + \sigma^2\mathbf{I}|^{\frac{1}{2}} |\Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T + \sigma^2\mathbf{I}|^{\frac{1}{2}}} \right)^{-\frac{1}{2}}.$$

With probability one, the matrices  $\Phi\mathbf{U}_i\mathbf{U}_i^T\Phi^T$  and  $\Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T$  have rank  $k$ , and the matrix  $(\Phi\mathbf{U}_i\mathbf{U}_i^T + \Phi\mathbf{U}_j\mathbf{U}_j^T)/2$  has rank  $\min\{M, 2k\}$ . Let  $\lambda_{i_l}$  and  $\lambda_{j_l}$  denote the nonzero eigenvalues of  $\Phi\mathbf{U}_i\mathbf{U}_i^T$  and  $\Phi\mathbf{U}_j\mathbf{U}_j^T$ , respectively, and let  $\lambda_{ij_l}$  denote the nonzero eigenvalues of the latter matrix. Then, we can write the pairwise bound as

$$\begin{aligned} P_e(\mathbf{U}_i, \mathbf{U}_j) &\leq \frac{1}{2} \left( \frac{(\sigma^2)^{M - \min\{M, 2k\}} \prod_{l=1}^{\min(2k, M)} (\lambda_{ij_l} + \sigma^2)}{\sqrt{(\sigma^2)^{M-k} \prod_{l=1}^k (\lambda_{i_l} + \sigma^2) \cdot (\sigma^2)^{M-k} \prod_{l=1}^k (\lambda_{j_l} + \sigma^2)}} \right)^{-\frac{1}{2}} \\ &= \frac{1}{2} \cdot \left( \frac{1}{\sigma^2} \right)^{-\frac{\min\{M-k, k\}}{2}} \cdot \left( \frac{\prod_{l=1}^{\min(2k, M)} (\lambda_{ij_l} + \sigma^2)}{\sqrt{\prod_{l=1}^k (\lambda_{i_l} + \sigma^2) \cdot \prod_{l=1}^k (\lambda_{j_l} + \sigma^2)}} \right)^{-\frac{1}{2}}. \end{aligned}$$

By construction,

$$\Phi\mathbf{U}_i\mathbf{U}_i^T\Phi^T + \Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T \geq \Phi\mathbf{U}_i\mathbf{U}_i^T\Phi^T, \Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T.$$

By Weyl's monotonicity theorem (see, e.g., [56]),  $2\lambda_{ij_l} \geq \lambda_{i_l}$  and  $2\lambda_{ij_l} \geq \lambda_{j_l}$  for every  $1 \leq l \leq k$ . Therefore,

$$\prod_{l=1}^k 2(\lambda_{ij_l} + \sigma^2) \geq \sqrt{\prod_{l=1}^k (\lambda_{i_l} + \sigma^2) \cdot \prod_{l=1}^k (\lambda_{j_l} + \sigma^2)},$$

from which it follows that

$$\begin{aligned} P_e(\mathbf{U}_i, \mathbf{U}_j) &\leq \frac{1}{2} \cdot \left(\frac{1}{\sigma^2}\right)^{-\frac{\min\{M-k, k\}}{2}} \cdot 2^{\frac{k}{2}} \cdot \left(\prod_{l=k+1}^{\min(2k, M)} (\lambda_{ij_l} + \sigma^2)\right)^{-\frac{1}{2}} \\ &\leq \frac{1}{2} \cdot \left(\frac{1}{\sigma^2}\right)^{-\frac{\min\{M-k, k\}}{2}} \cdot 2^{\frac{k}{2}} \cdot (\lambda_{ij_{\min\{2k, M\}}} + \sigma^2)^{-\frac{\min\{M-k, k\}}{2}} \\ &= 2^{\frac{k-2}{2}} \cdot \left(1 + \frac{\lambda_{ij_{\min\{2k, M\}}}}{\sigma^2}\right)^{-\frac{\min\{M-k, k\}}{2}}. \end{aligned}$$

Next, we bound the eigenvalue  $\lambda_{ij_{\min\{2k, M\}}}$ . Because each matrix  $\mathbf{U}_i$  has i.i.d. Gaussian entries with zero mean and variance  $1/k$ , so too does each matrix  $\Phi\mathbf{U}_l \in \mathbb{R}^{M \times k}$ . Furthermore, observe that

$$\Phi\mathbf{U}_i\mathbf{U}_i^T\Phi^T + \Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T = [\Phi\mathbf{U}_i \quad \Phi\mathbf{U}_j] \cdot \begin{bmatrix} (\Phi\mathbf{U}_i)^T \\ (\Phi\mathbf{U}_j)^T \end{bmatrix}.$$

Therefore, the nonzero eigenvalues of  $\Phi\mathbf{U}_i\mathbf{U}_i^T\Phi^T + \Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T$  are those of a scaled Wishart matrix. Specifically, if  $2k < M$ , the eigenvalues are those of  $1/k \cdot \mathbf{W}_1$ , where  $\mathbf{W}_1 \sim \mathcal{W}_{2k}(M, \mathbf{I})$ . If  $M \geq 2k$ , the eigenvalues are those of  $1/k \cdot \mathbf{W}_2$ , where  $\mathbf{W}_2 \sim \mathcal{W}_M(2k, \mathbf{I})$ . By [55, Theorem 1], the minimum eigenvalue in either case converges on  $(\sqrt{1/\kappa} - \sqrt{2})^2$  almost surely. Therefore,  $\lambda_{ij_{\min\{2k, M\}}}$  converges on  $(\sqrt{1/(2\kappa)} - 1)^2$  almost surely, and we obtain

$$P_e(\mathbf{U}_i, \mathbf{U}_j) \leq 2^{\frac{k-2}{2}} \cdot \left(1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2}\right)^{-\frac{\min\{M-k, k\}}{2}}, \quad (26)$$

where  $\epsilon(M) \rightarrow 0$  almost surely as  $M \rightarrow \infty$ . *A fortiori*, the bound in (26) is also a bound on the *expected* pairwise probability, with  $\epsilon(M)$  independent for each  $i, j$  pair.

Invoking the union bound over all  $L(M)$  subspaces, we obtain

$$\begin{aligned} E[P_e(\mathbf{a})] &\leq \frac{1}{L(M)} \sum_{l=1}^{L(M)} \sum_{l' \neq l} E[P_e(\mathbf{U}_l, \mathbf{U}_{l'})] \\ &= (L(M) - 1)E[P_e(\mathbf{U}_l, \mathbf{U}_{l'})] \\ &\leq 2^{\rho M} E[P_e(\mathbf{U}_l, \mathbf{U}_{l'})], \end{aligned}$$

where the second equality follows because each  $\mathbf{U}_l$  is drawn independently. Taking the logarithm of both sides yields

$$\log_2(E[P_e(\mathbf{a})]) \leq \rho M + \frac{k-2}{2} - \frac{\min\{M-k, k\}}{2} \log_2 \left(1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2}\right). \quad (27)$$

Therefore, if

$$\rho < \frac{\min\{1 - \kappa, \kappa\}}{2} \log_2 \left(1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2}{\sigma^2}\right) - \frac{\kappa}{2},$$

then  $E[P_e(\mathbf{a})]$  goes to zero as  $M \rightarrow \infty$ , and thus  $P_e(\mathbf{a})$  goes to zero almost surely, as was to be shown.  $\blacksquare$

When  $\kappa \geq \frac{1}{2}$ , the lower and upper bounds agree to within a  $O(1)$  term; otherwise they are loose. Based on the numerical experiments presented in Section V, we conjecture that the upper bound is approximately tight, while the lower bound is loose.

## B. Affine Subspaces

Next, we bound  $C_{\text{affine}}(\nu, \kappa)$ .

*Theorem 2:* For affine subspace classification, the classification capacity satisfies

$$C_{\text{affine}}(\nu, \kappa) \leq \frac{1-\kappa}{2} \log_2 \left(\frac{1}{\sigma^2}\right) + \frac{1}{2} \log_2(2 + \sigma^2) - \frac{\kappa}{2} \log_2((\sqrt{1/\kappa} - 1)^2 + \sigma^2), \quad (28)$$

and

$$C_{\text{affine}}(\nu, \kappa) \geq \begin{cases} \frac{1-\kappa}{2} \log_2 \left( 1 + \frac{\min\{(\sqrt{1/(2\kappa)}-1)^2, 1/2\}}{\sigma^2} \right) - \frac{\kappa}{2} & \text{for } \kappa < 1/2 \\ \frac{1-\kappa}{2} \log_2 \left( 1 + \frac{(\sqrt{1/(2\kappa)}-1)^2}{\sigma^2} \right) - \frac{\kappa}{2} & \text{for } \kappa \geq 1/2 \end{cases}. \quad (29)$$

*Proof:* As before, we prove the upper bound by bounding the mutual information, and the lower bound by direct analysis of the probability of error via the Bhattacharyya bound.

**Upper Bound:** To prove the upper bound, we expand the mutual information as

$$I(\mathbf{a}; \mathbf{y}) = I(\mathbf{U}, \mu; \mathbf{y}) = h(\mathbf{y}) - h(\mathbf{y}|\mathbf{U}, \mu). \quad (30)$$

As in the case of linear subspaces, the  $\Phi$  that maximizes the mutual information has unit singular values. Furthermore, because the entropy of a Gaussian does not depend on the mean,  $h(\mathbf{y}|\mathbf{U}, \mu) = h(\mathbf{y}|\mathbf{U})$ . Therefore, applying (23), we obtain

$$h(\mathbf{y}|\mathbf{U}, \mu) \geq \frac{k}{2} \log_2((\sqrt{1/\kappa} - 1)^2 + \epsilon(M) + \sigma^2) + \frac{M-k}{2} \log_2(\sigma^2) + \frac{M}{2} \log_2(2\pi e), \quad (31)$$

where  $\epsilon(M) \rightarrow 0$ . Then, observing that

$$\begin{aligned} E[\mathbf{y}\mathbf{y}^T] &= E_{\mu}[\Phi\mu\mu^T\Phi^T] + E_{\mathbf{U}}[\Phi\mathbf{U}\mathbf{U}^T\Phi^T] + \sigma^2\mathbf{I} \\ &= (2 + \sigma^2)\mathbf{I}, \end{aligned}$$

we conclude that

$$h(\mathbf{y}) \leq \frac{M}{2} \log_2(2\pi e(2 + \sigma^2)), \quad (32)$$

from which it follows that

$$\lim_{M \rightarrow \infty} \frac{I(\mathbf{y}; \mathbf{U}, \mu)}{M} \leq \frac{1-\kappa}{2} \log \left( \frac{1}{\sigma^2} \right) + \frac{1}{2} \log_2(2 + \sigma^2) - \frac{\kappa}{2} \log_2((\sqrt{1/\kappa} - 1)^2 + \sigma^2). \quad (33)$$

Applying the preceding to Lemma 1, we obtain the upper bound.

**Lower Bound:** Suppose that we choose  $\Phi$  to be any matrix with orthonormal rows. We bound the pairwise misclassification error via Lemma 2, which yields

$$\begin{aligned} P_e(\mu_i, \mathbf{U}_i, \mu_j, \mathbf{U}_j) &\leq \frac{1}{2} \cdot \left( \frac{|\frac{\Phi\mathbf{U}_i\mathbf{U}_i^T + \Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T + 2\sigma^2\mathbf{I}}{2}|}{|\Phi\mathbf{U}_i\mathbf{U}_i^T\Phi + \sigma^2\mathbf{I}|^{\frac{1}{2}} |\Phi\mathbf{U}_j\mathbf{U}_j^T\Phi^T + \sigma^2\mathbf{I}|^{\frac{1}{2}}} \right)^{-\frac{1}{2}} \\ &\quad \exp \left( -\frac{1}{8} \cdot (\mu_i - \mu_j)^T \Phi^T \left( \frac{\Phi(\mathbf{U}_i\mathbf{U}_i^T + \mathbf{U}_j\mathbf{U}_j^T)\Phi^T + 2\sigma^2\mathbf{I}}{2} \right)^{-1} \Phi(\mu_i - \mu_j) \right). \quad (34) \end{aligned}$$

Observe that the argument of the exponential term is always nonnegative, so the exponential is always smaller than one. Therefore, the bound on the misclassification probability of affine subspaces is always smaller than that of linear subspaces, and the lower bound from Theorem 1 also applies to affine subspaces. Applying this fact yields the lower bound for  $\kappa \geq 1/2$ .

For  $\kappa < 1/2$ , we apply (26) to (34), yielding

$$\begin{aligned} P_e(\mu_i, \mathbf{U}_i, \mu_j, \mathbf{U}_j) &\leq 2^{\frac{\kappa-2}{2}} \cdot \left( 1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2} \right)^{-\frac{\kappa}{2}} \\ &\quad \exp \left( -\frac{1}{8} \cdot (\mu_i - \mu_j)^T \Phi^T \left( \frac{\Phi(\mathbf{U}_i\mathbf{U}_j^T + \mathbf{U}_i\mathbf{U}_j^T)\Phi^T + 2\sigma^2\mathbf{I}}{2} \right)^{-1} \Phi(\mu_i - \mu_j) \right), \quad (35) \end{aligned}$$

where again  $\epsilon(M) \rightarrow 0$ . Next, let  $\Phi(\mathbf{U}_i\mathbf{U}_i^T + \mathbf{U}_j\mathbf{U}_j^T)\Phi^T = \mathbf{W}_{ij}\Lambda_{ij}\mathbf{W}_{ij}^T$  be the eigenvalue decomposition of the covariance pair sum. Also define  $\omega = \mathbf{W}_{ij}^T\Phi(\mu_i - \mu_j)/2$ , which is i.i.d. Gaussian with zero mean and unit variance.

We therefore obtain

$$P_e(\mu_i, \mathbf{U}_i, \mu_j, \mathbf{U}_j) \leq 2^{\frac{k-2}{2}} \cdot \left( 1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2} \right)^{-\frac{k}{2}} \cdot \exp\left(-\frac{1}{4}\omega^T (\Lambda_{ij}/2 + \sigma^2 \mathbf{I})^{-1} \omega\right).$$

With probability one,  $\Lambda_{ij}$  contains  $2k$  nonzero eigenvalues. The preceding bound increases in these eigenvalues, so to bound the error we bound the eigenvalues by infinity, which yields

$$P_e(\mu_i, \mathbf{U}_i, \mu_j, \mathbf{U}_j) \leq 2^{\frac{k-2}{2}} \cdot \left( 1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2} \right)^{-\frac{k}{2}} \cdot \prod_{i=2k+1}^M \exp\left(-\frac{1}{4\sigma^2}\omega_i^2\right).$$

Taking the expectation yields

$$E[P_e(\mu_i, \mathbf{U}_i, \mu_j, \mathbf{U}_j)] \leq 2^{\frac{k-2}{2}} \cdot \left( 1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2} \right)^{-\frac{k}{2}} \cdot \prod_{i=2k+1}^M E\left[\exp\left(-\frac{1}{4\sigma^2}\omega_i^2\right)\right],$$

where the expectation moves inside the product because each  $\omega_i$  is independent of the others. Noting that each expectation in the final expression is just the moment-generating function of a Chi-squared random variable, we obtain

$$E[P_e(\mu_i, \mathbf{U}_i, \mu_j, \mathbf{U}_j)] \leq 2^{\frac{k-2}{2}} \cdot \left( 1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2} \right)^{-\frac{k}{2}} \cdot \left( 1 + \frac{1}{2\sigma^2} \right)^{-\frac{M-2k}{2}} \quad (36)$$

Applying the union bound and taking the logarithm, we obtain

$$\log_2(E[P_e(\mathbf{a})]) \leq \rho M - \frac{k}{2} \log_2 \left( 1 + \frac{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M)}{\sigma^2} \right) - \frac{M-2k}{2} \log_2 \left( 1 + \frac{1}{2\sigma^2} \right) + \frac{k-2}{2} \quad (37)$$

$$\leq \rho M - \frac{M-k}{2} \log_2 \left( 1 + \frac{\min\{(\sqrt{1/(2\kappa)} - 1)^2 + \epsilon(M), 1/2\}}{\sigma^2} \right) + \frac{k-2}{2}. \quad (38)$$

Letting  $M \rightarrow \infty$ , we obtain the lower bound for  $\kappa < 1/2$ . ■

For affine subspaces, the bounds are tight to within an  $O(1)$  term for all values of  $\kappa$ . Roughly speaking, the term in the Bhattacharyya bound associated with discriminating the means cancels out the gap to the upper bound associated with discriminating the associated linear subspaces. Therefore pairwise analysis, along with the union bound, is sufficient for establishing tight bounds on the classification capacity for affine subspaces even when it fails for linear subspaces.

#### IV. DIVERSITY-DISCRIMINATION TRADEOFF

Here we prove bounds on the diversity-discrimination tradeoff. Similar to the classification capacity, we prove an upper bound on  $d_{\text{linear}}(r)$ , which shows that the maximum diversity gain is  $\min\{k, M-k\}$  and the maximum discrimination gain is  $M-k$ . We also prove a lower bound, based on the Bhattacharyya bound, which establishes that the average misclassification probability decays at least as  $(1/\sigma^2)^{-\frac{\min\{k, M-k\}-r}{2}}$  when the number of subspaces grows as  $(1/\sigma^2)^{\frac{r}{2}}$ . For affine subspaces, we prove an upper bound which shows that the misclassification probability decays no faster than  $(1/\sigma^2)^{-\frac{M-k-r}{2}}$ . In this case, the Bhattacharyya analysis shows that the upper bound is tight.

##### A. Linear Subspaces

First, we prove bounds on  $d_{\text{linear}}(r)$ .

*Theorem 3:* For linear subspaces, the DDT is upper bounded by

$$d_{\text{linear}}(r) \leq \left[ \min \left\{ M - k - r, k \left( 1 - \frac{r}{M} \right) \right\} \right]^+, \quad (39)$$

and the DDT is bounded below by

$$d_{\text{linear}}(r) \geq [\min\{M-k, k\} - r]^+. \quad (40)$$

*Proof:* First we prove the upper bound, the first term of which follows from Lemma 3. Combining (23) and (24), it is easy to see that

$$\lim_{\sigma^2 \rightarrow 0} \frac{I(\mathbf{U}; \mathbf{y})}{\frac{1}{2} \log_2(1/\sigma^2)} \leq M - k. \quad (41)$$

Therefore, by Lemma 3,  $d_{\text{linear}}(r) = 0$  whenever  $r \geq M - k$ . Next, suppose that  $d_{\text{linear}}(r) > M - k - r$  for some  $0 \leq r < M - k$ , meaning that, for some  $\epsilon > 0$ ,

$$\log_2(\bar{P}_\epsilon(r, \sigma^2)) \leq -\frac{M - k - r + \epsilon}{2} \log_2(1/\sigma^2) + o(\log(\sigma^2)). \quad (42)$$

Using the union bound, we can express the probability of error for  $r = M - k$  in terms of (42):

$$\begin{aligned} \log_2(\bar{P}_\epsilon(M - k, \sigma^2)) &\leq \log_2((1/\sigma^2)^{\frac{M-k-r}{2}} \bar{P}_\epsilon(r, \sigma^2)) \\ &\leq \frac{M - k - r}{2} \log_2(1/\sigma^2) + \log_2(\bar{P}_\epsilon(r, \sigma^2)) \\ &\leq -\frac{\epsilon}{2} \log_2(1/\sigma^2) + o(\log(\sigma^2)). \end{aligned}$$

This implies  $d_{\text{linear}}(M - k) > 0$ , which is a contradiction.

The second term in the upper bound follows from an ‘‘outage’’-style argument reminiscent of that of [14]. For linear subspaces, we can rewrite the signal model (1) as

$$\mathbf{y} = \Phi \mathbf{U} \mathbf{h} + \mathbf{z},$$

where  $\mathbf{h} = (h_1, \dots, h_k)^T \sim \mathcal{N}(0, \mathbf{I})$ . We define an outage event

$$\mathcal{F} = \{h_i^2 \leq (1/\sigma^2)^{-\beta}, \forall i\}, \quad (43)$$

for  $0 \leq \beta \leq 1$ . Because each  $h_i^2$  is Chi squared with a single degree of freedom,

$$\Pr(\mathcal{F}) \leq (1/\sigma^2)^{-\frac{k\beta}{2}} \cdot \exp(k/2). \quad (44)$$

Next, we bound the *conditional* normalized mutual information:

$$\begin{aligned} \lim_{\sigma^2 \rightarrow 0} \frac{I(\mathbf{U}; \mathbf{y} | \mathcal{F})}{1/2 \log_2(1/\sigma^2)} &= \lim_{\sigma^2 \rightarrow 0} \frac{h(\mathbf{y} | \mathcal{F}) - h(\mathbf{y} | \mathbf{U}, \mathcal{F})}{1/2 \log_2(1/\sigma^2)} \\ &\leq \lim_{\sigma^2 \rightarrow 0} \frac{\log_2(E_{\mathbf{h}, \mathbf{U}}[\det(\Phi \mathbf{U} \mathbf{h} \mathbf{h}^T \mathbf{U}^T \Phi^T + \sigma^2 \mathbf{I}) | \mathcal{F}]) - M/2 \log_2(\sigma^2)}{1/2 \log_2(1/\sigma^2)}, \end{aligned}$$

where the inequality follows because (i) the Gaussian distribution maximizes mutual information, and (ii)  $h(\mathbf{y} | \mathbf{U}, \mathcal{F}) \geq h(\mathbf{z})$  by the entropy power inequality. Conditioned on the outage event, we have  $\mathbf{h} \mathbf{h}^T \leq (1/\sigma^2)^{-\beta} \cdot \mathbf{I}$ , from which it follows that

$$\begin{aligned} \lim_{\sigma^2 \rightarrow 0} \frac{I(\mathbf{U}; \mathbf{y} | \mathcal{F})}{1/2 \log_2(1/\sigma^2)} &\leq \lim_{\sigma^2 \rightarrow 0} \frac{\log_2(E_{\mathbf{U}}[|(1/\sigma^2)^{-\beta} \cdot \Phi \mathbf{U} \mathbf{U}^T \Phi^T + \sigma^2 \mathbf{I}|]) - M/2 \log_2(\sigma^2)}{1/2 \log_2(1/\sigma^2)} \\ &\leq \lim_{\sigma^2 \rightarrow 0} \frac{M/2 \log_2((1/\sigma^2)^{-\beta} + \sigma^2) - M/2 \log_2(\sigma^2)}{1/2 \log_2(1/\sigma^2)} \\ &= M(1 - \beta). \end{aligned}$$

By the law of total probability,

$$\begin{aligned} E[\mathbb{P}_e(\mathbf{a})] &= E[\mathbb{P}_e(\mathbf{a}) | \mathcal{F}] \Pr(\mathcal{F}) + E[\mathbb{P}_e(\mathbf{a}) | \mathcal{F}^c] (1 - \Pr(\mathcal{F})) \\ &\geq E[\mathbb{P}_e(\mathbf{a}) | \mathcal{F}] \Pr(\mathcal{F}). \end{aligned}$$

By Lemma 3, whenever  $r > M(1 - \beta)$ , the conditional probability  $E[P(\mathbf{a})|\mathcal{F}]$  is bounded away from zero. Therefore,

$$\begin{aligned} d_{\text{linear}}(M(1 - \beta) + \epsilon) &= \lim_{\sigma^2 \rightarrow 0} -\frac{\log_2(E[P_e(\mathbf{a})])}{1/2 \log_2(1/\sigma^2)} \\ &\leq \lim_{\sigma^2 \rightarrow 0} -\frac{\log_2(E[P(\mathbf{a})|\mathcal{F}]) + \log_2(\Pr(\mathcal{F}))}{1/2 \log_2(1/\sigma^2)} \\ &= \lim_{\sigma^2 \rightarrow 0} -\frac{\log_2(\Pr(\mathcal{F}))}{1/2 \log_2(1/\sigma^2)} \\ &\leq k\beta. \end{aligned}$$

Taking  $\epsilon \rightarrow 0$ , we obtain the second term of the upper bound.

The lower bound follows from the Bhattachayya analysis from Section III-A. For discrimination gain  $r$ , we combine (26) with the union bound, yielding

$$\log_2(\bar{P}_e(r, \sigma^2)) \leq \frac{r}{2} \log_2\left(\frac{1}{\sigma^2}\right) - \frac{\min\{M - k, k\}}{2} \log_2\left(1 + \frac{(\sqrt{1/\kappa} - \sqrt{2})^2 + \epsilon(M)}{\sigma^2}\right) + \frac{k}{2}.$$

Therefore, we have

$$d_{\text{linear}}(r) = \lim_{\sigma^2 \rightarrow 0} \frac{\log_2(\bar{P}_e(r, \sigma^2))}{\frac{1}{2} \log_2(1/\sigma^2)} \geq \min\{M - k, k\} - r. \quad \blacksquare$$

Similar to the classification capacity, the lower bound is tight when  $k \geq M/2$ . Otherwise, the lower bound achieves full diversity for  $r = 0$ , but falls short of the upper bound for higher discrimination gain. Note, however, that the second term in (39), which establishes that the diversity gain is no greater than  $k$ , is clearly loose because it predicts nonzero diversity for discrimination gains higher than  $M - k$ . This looseness is due to the bound on the normalized mutual information, in which we employed  $h(\mathbf{y}|\mathbf{U}, \mathcal{F}) \geq h(\mathbf{z})$ ; this bound neglects the effect of the outage event on the eigenvalues.

A tighter bound on the conditional entropy is difficult because  $\mathbf{y}$  is no longer Gaussian conditioned on  $\mathcal{F}$ . However, we can make heuristic calculations by bounding the conditional covariance and supposing that the entropy is approximately that of the equivalent Gaussian. Then, the normalized mutual information is instead bounded by  $(M - k)(1 - \beta)$ . Following that analysis leads to the following bound on DDT function

$$d_{\text{linear}}(r) = \min\{M - k, k\} \left[1 - \frac{r}{M - k}\right]^+. \quad (45)$$

This function is just the line segment connecting the maximum diversity order and the maximum discrimination gain in the upper bound. Based on the preceding intuition and the numerical results in Section V, we conjecture that this is the true diversity-discrimination tradeoff for linear classification.

### B. Affine Subspaces

Next, we derive  $d_{\text{affine}}(r)$ .

*Theorem 4:* For affine subspace classification, the DDT is

$$d_{\text{affine}}(r) = [M - k - r]^+. \quad (46)$$

*Proof:* We can upper bound the DDT using the same argument as in the proof of Theorem 3. Combining (31) and (32), we obtain

$$\lim_{\sigma^2 \rightarrow 0} \frac{I(\mu, \mathbf{U}; \mathbf{y})}{\frac{1}{2} \log_2(1/\sigma^2)} \leq M - k. \quad (47)$$

Therefore,  $d_{\text{affine}}(r) = 0$  for  $r \geq M - k$  by Lemma 3. As before, by the union bound there is a contradiction if  $d(r) > M - k - r$  for any  $0 \leq r \leq M - k$ .

To lower bound the DDT, observe from the proof of Theorem 2 that

$$\log_2(E[\mathbb{P}_e(\mu_i, \Sigma_i, \mu_j, \Sigma_j)]) \leq -\frac{M-k}{2} \log_2(1/\sigma^2) + o(\log(\sigma^2)). \quad (48)$$

For discrimination gain  $0 \leq r \leq M-k$ , the union bound yields

$$\log_2(\bar{\mathbb{P}}_e(r, \sigma^2)) \leq -\frac{M-k-r}{2} \log_2(1/\sigma^2) + o(\log(\sigma^2)), \quad (49)$$

which establishes the result.  $\blacksquare$

For affine spaces, similar to the classification capacity, the upper bound is tight. Therefore, the Bhattacharyya bound is tight not only with respect to the pairwise error, but also upon application of the union bound. Therefore, while the translation of linear subspaces into affine subspaces does not necessarily improve the *number* of subspaces that a classifier can discriminate reliably, for  $k < M/2$  translation does improve the decay of the error probability as the noise power vanishes.

## V. NUMERICAL RESULTS

Here we validate our results numerically. First, we study the performance of linear subspace classifiers with subspaces drawn randomly from  $p_{\text{linear}}$ , focusing on the regimes in which the upper and lower bounds disagree and drawing conclusions about the tightness of the bounds. Then, we study the classifier performance over the YaleB database of face images, comparing empirical performance to the predictions of Section III.

### A. Linear Subspaces

Here, we want to see whether the upper bound on the  $C_{\text{linear}}(\nu, \kappa)$  is tight for  $\kappa < 1/2$ , as we conjectured in Section III-A. We also want to see whether the diversity-discrimination function we conjectured in Section IV-A is correct.

To answer these questions, we examine classifier performance as  $\sigma^2 \rightarrow 0$ . We draw subspaces from  $p_{\text{linear}}$  and we choose  $\Phi$  to be the first  $M$  rows of a randomly-chosen unitary matrix. We corrupt the features with white Gaussian noise of variance  $\sigma^2$ , and we perform maximum-likelihood classification on the noisy features. Because of computational limitations, it is infeasible to study empirical performance as the signal dimension becomes large. Therefore, instead of testing the classification capacity directly, we examine the DDT performance. If, for discrimination gain  $r$ , the diversity gain is nonzero, then the classification capacity must be at least as great as  $r/M \log_2(1/\sigma^2) + o(\log(\sigma^2))$ .

In Figure 1 we plot the misclassification probability as a function of  $\sigma^2$ . For each value of  $\sigma^2$ , we compute the average misclassification probability over  $10^2$  realizations of the subspaces  $\mathbf{U}_l$  and  $10^2$  realizations of the signal of interest per set of subspaces. We also plot the error slopes predicted by (45). We select dimensions  $N = M = 3$ ,  $k = 1$  and discrimination gains  $r \in \{0, 0.75, 1.5, 1.8\}$ . We observe decaying misclassification probability for all values of  $r$ ; furthermore, we observe rates of decay consistent with the conjectured DDT function. Therefore, we conclude that, regardless of  $\kappa$ , the classification capacity satisfies

$$C_{\text{linear}}(\nu, \kappa) = \frac{1-\kappa}{2} \log_2(1/\sigma^2) + o(\log_2(\sigma^2)),$$

and that the DDT is

$$d_{\text{linear}}(r) = \min\{k, M-k\} \left[ 1 - \frac{r}{M-k} \right]^+.$$

Recall that these conclusions were proven only in the regimes  $k \geq M/2$  or  $\kappa \geq 1/2$ .

### B. Face Recognition

Next, we explore the correspondence between the theoretical results derived in the previous sections and a practical face recognition application. We examine face recognition when the orientation of the face relative to the camera remains fixed but the illumination varies. Supposing the faces to be approximately convex and to reflect light according to Lambert's law, [57] shows via spherical harmonics that the set of images of an individual face lies

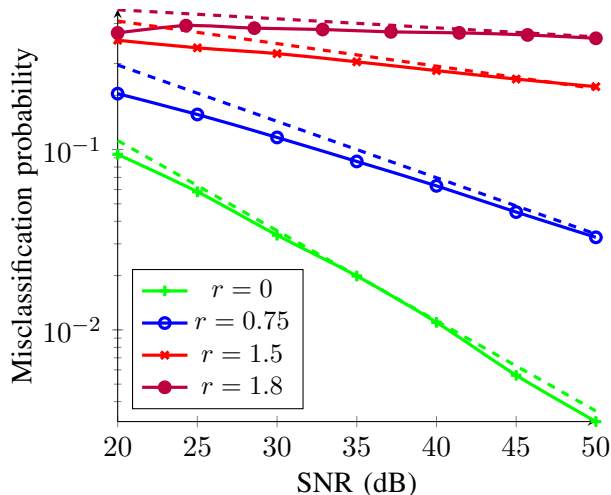


Fig. 1. Misclassification probability vs. SNR for linear subspaces. Dashed lines indicate the slopes predicted by (45).

approximately on a nine-dimensional subspace, regardless of the inherent dimension of the images. It is therefore sufficient to discriminate between the subspaces to classify faces.

We use 38 cropped faces from the Extended Yale Face Database B, described in [7], [8]. For each face, the database contains a few dozen greyscale photographs, each having 32,256 pixels, taken under a variety of illumination conditions as shown in Figure 2. We vectorize these images and pass them through a feature matrix  $\Phi$ , chosen as before to be the first  $M$  rows of an arbitrary unitary matrix. We classify the faces using the maximum-likelihood classifier supposing zero-mean Gaussian classes. We divide the database into two, using half of the images to estimate the nearest covariance for each face, using the other half as test images.

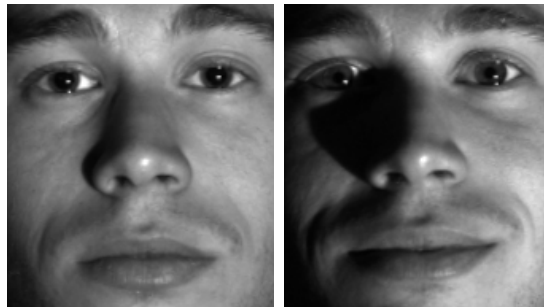


Fig. 2. Two sample images from the Extended Yale Face Database B. These images are of the same face, but are taken under different illumination conditions.

In Figure 3 we plot the misclassification probability as a function of  $M$  and for  $L$  ranging from 2 to 38. While we do not label each curve, it is easy to see that the misclassification probability increases with  $L$  and decreases with  $M$ . However, even for large  $M$  the error probability remains as high as 0.2 for  $L = 38$ . We take 0.2 as a baseline for “successful” performance when the number of faces and signal dimension are high.

Finally, we examine how well our theory predicts the performance seen here. To estimate the noise power  $\sigma^2$ , we project each image onto its estimated subspace, transformed by  $\Phi$ , and we take the projected squared norm as the signal power and the squared residual norm, normalized by the number of features  $M$ , as the noise power. We then estimate the number of faces that Theorem 1 predicts can be discriminated reliably. Discarding the constants, we simply compute

$$\max\{1, \min\{1/\sigma^{2(M-9)/2}, 38\}\}. \quad (50)$$



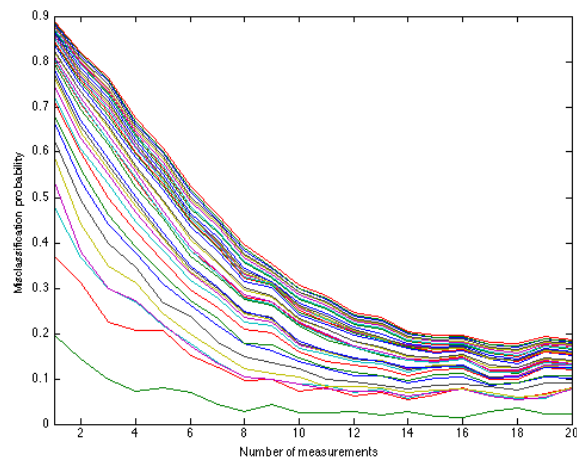


Fig. 3. Misclassification probability as a function of  $M$ , for  $L$  ranging from 1 to 38.

Naturally, this number grows quickly in  $M$ , and beyond  $M = 11$  or  $M = 12$ , theory suggests that we ought to be able to discriminate all 38 of the faces with low probability of error. In Figure 4 we compare this prediction against the empirical performance of our classifier. Using the results shown in Figure 3, we compute, for each  $M$ , the maximum  $L$  for which the probability of error is less than 0.2.

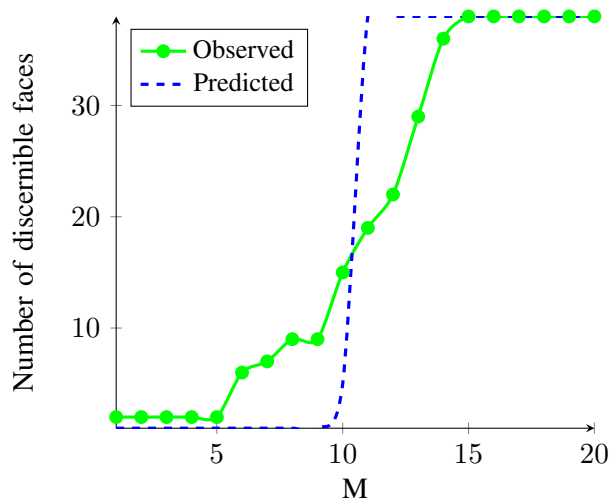


Fig. 4. Misclassification probability as a function of  $M$ , for  $L$  ranging from 1 to 38.

The empirical performance is similar to theoretical prediction. As  $M$  increases past 9, the number of faces rises swiftly as predicted. After  $M = 15$  or so, all 38 of the faces can be discriminated, and it is not advantageous to extract more features. We do observe, however, that the transition is not as sharp as Theorem 1 predicts. Whereas the theoretical transition occurs over only 2-3 features, in practice the transition stretches out over 5-10 features. In addition to mild model mismatch due to non-Lambertian reflectances, shadows due to the non-convexity of real faces, imperfect estimation of subspaces, etc., we suspect that this is primarily a phenomenon of classification at finite dimension. The transition between success and failure becomes sharp in the limit, but remains gradual when dimensions measure in the tens or hundreds.

## VI. CONCLUSION

Inspired by dualities between wireless communication over non-coherent channels and the classification from noisy, linear features, we have derived fundamental limits on the classification of linear and affine subspaces from noisy, linear features. We defined performance limits reminiscent of those in wireless information theory: the classification capacity, which governs classifier performance in the limit of high signal dimension, and the diversity-discrimination tradeoff, which governs classifier performance in the limit of low noise power. We proved inner and outer bounds on these quantities. For linear subspaces, the bounds are tight in some regimes of  $N$ ,  $M$ , and  $k$ , and for affine subspaces they are tight everywhere. Based on numerical evaluation, we conjectured that the true classification capacity and DDT for linear subspaces in the regimes in which the bounds are not tight. Beyond the characterization of such limits, we showed via an application to face recognition that theoretical trends agree reasonably with practical ones.

## REFERENCES

- [1] Y. Adini, Y. Moses, and S. Ullman, "Face recognition: The problem of compensating for changes in illumination direction," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 19, no. 7, pp. 721–732, 1997.
- [2] J. J. Hull, "A database for handwritten text recognition research," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 16, no. 5, pp. 550–554, 1994.
- [3] Y. LeCun, L. Jackel, L. Bottou, A. Brunot, C. Cortes, J. Denker, H. Drucker, I. Guyon, U. Muller, E. Sackinger *et al.*, "Comparison of learning algorithms for handwritten digit recognition," in *International Conference on Artificial Neural Networks*, vol. 60, 1995.
- [4] D. T. Ross, U. Scherf, M. B. Eisen, C. M. Perou, C. Rees, P. Spellman, V. Iyer, S. S. Jeffrey, M. Van de Rijn, M. Waltham *et al.*, "Systematic variation in gene expression patterns in human cancer cell lines," *Nature Genetics*, vol. 24, no. 3, pp. 227–235, 2000.
- [5] U. Alon, N. Barkai, D. A. Notterman, K. Gish, S. Ybarra, D. Mack, and A. J. Levine, "Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays," *Proceedings of the National Academy of Sciences*, vol. 96, no. 12, pp. 6745–6750, 1999.
- [6] [Online]. Available: <http://www.shazam.com/>
- [7] A. S. Georghiades, P. N. Belhumeur, and D. J. Kriegman, "From few to many: Illumination cone models for face recognition under variable lighting and pose," *IEEE trans. Pattern Analysis and Machine Intelligence*, vol. 23, no. 6, pp. 643–660, 2001.
- [8] K.-C. Lee, J. Ho, and D. J. Kriegman, "Acquiring linear subspaces for face recognition under variable lighting," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 27, no. 5, pp. 684–698, 2005.
- [9] R. Epstein, P. W. Hallinan, and A. L. Yuille, "5±2 eigenimages suffice: An empirical investigation of low-dimensional lighting models," in *Proc. the Workshop on Physics-Based Modeling in Computer Vision, 1995*. IEEE, 1995.
- [10] R. A. Fisher, "The use of multiple measurements in taxonomic problems," *Annals of Eugenics*, vol. 7, no. 2, pp. 179–188, 1936.
- [11] L. Zheng and D. Tse, "Diversity and multiplexing: a fundamental tradeoff in multiple-antenna channels," *IEEE Trans. Inform. Theory*, vol. 49, no. 5, pp. 1073–1096, May 2003.
- [12] L. Zheng and D. N. C. Tse, "Communication on the Grassmann manifold: A geometric approach to the noncoherent multiple-antenna channel," *IEEE Trans. Info. Theory*, vol. 48, no. 2, pp. 359–383, 2002.
- [13] T. L. Marzetta and B. M. Hochwald, "Capacity of a mobile multiple-antenna communication link in Rayleigh flat fading," *IEEE Trans. Info. Theory*, vol. 45, no. 1, pp. 139–157, 1999.
- [14] L. Zheng and D. N. C. Tse, "The diversity-multiplexing tradeoff for non-coherent multiple antenna channels," in *Proc. Allerton*, vol. 40, no. 2, 2002, pp. 834–843.
- [15] W. Yang, G. Durisi, and E. Riegler, "On the capacity of large-MIMO block-fading channels," *IEEE J. Select Areas Commun.*, vol. 31, no. 2, pp. 117–132, 2013.
- [16] M. Nokleby, M. Rodrigues, and R. Calderbank, "Information-theoretic limits on the classification of Gaussian mixtures: Classification on the Grassmann manifold," in *Proc. Information Theory Workshop (ITW)*, 2013.
- [17] R. Dobrushin, "General formulation of shannon's main theorem in information theory," *Trans. American Mathematical Society*, vol. 33, pp. 323–438, 1963.
- [18] M. Chen, W. Carson, M. R. D. Rodrigues, R. Calderbank, and L. Carin, "Communication inspired linear discriminant analysis," in *Proceedings of the 29th International Conference on Machine Learning*, June 2012.
- [19] D. Erdogmus and J. C. Principe, "Lower and upper bounds for misclassification probability based on renyi's information," *VLSI signal processing systems for signal, image and video technology, Journal of*, vol. 37, no. 2 -3, pp. 305–317, 2004.
- [20] K. Hild, D. Erdogmus, K. Torkkola, and J. Principe, "Feature extraction using information-theoretic learning," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 28, no. 9, pp. 1385–1392, 2006.
- [21] S. Kaski and J. Peltonen, "Informative discriminant analysis," in *Proceedings of the Twentieth International Conference on Machine Learning (ICML)*, 2003, pp. 329–336.
- [22] L. Liu and P. Fieguth, "Texture classification from random features," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 34, no. 3, pp. 574–586, March 2012.
- [23] Z. Nenadic, "Information discriminant analysis: Feature extraction with an information-theoretic objective," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 29, no. 8, pp. 1394–1407, August 2007.
- [24] D. Tao, X. Li, X. Wu, and S. Maybank, "Geometric mean for subspace selection," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 31, no. 2, pp. 260–274, 2009.
- [25] K. Torkkola, "Learning discriminative feature transforms to low dimensions in low dimensions," in *Advances in neural information processing systems 14 (NIPS)*. MIT Press, 2001, pp. 3–8.

- [26] —, “Feature extraction by non-parametric mutual information maximization,” *Machine Learning Research, Journal of*, vol. 3, pp. 1415–1438, March 2003.
- [27] Q. Qiu and G. Sapiro, “Learning robust subspace clustering,” *CoRR*, vol. abs/1308.0273, 2013.
- [28] J. Hamm and D. D. Lee, “Grassmann discriminant analysis: a unifying view on subspace-based learning,” in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 376–383.
- [29] J. Tenenbaum, V. de Silva, and J. Langford, “A global geometric framework for nonlinear dimensionality reduction,” *Science*, vol. 290, no. 5500, pp. 2319–2323, December 2000.
- [30] L. Wang, A. Razi, M. D. Rodrigues, and R. Calderbank, “Nonlinear information-theoretic compressive measurement design,” in *Proc. International Conference on Machine Learning*, 2014.
- [31] S. Aeron, V. Saligrama, and M. Zhao, “Information theoretic bounds for compressed sensing,” *Information Theory, IEEE Transactions on*, vol. 56, no. 10, pp. 5111–5130, 2010.
- [32] M. Akcakaya and V. Tarokh, “Shannon-theoretic limits on noisy compressive sampling,” *Information Theory, IEEE Transactions on*, vol. 56, no. 1, pp. 492–504, 2010.
- [33] A. Fletcher, S. Rangan, and V. Goyal, “Necessary and sufficient conditions for sparsity pattern recovery,” *Information Theory, IEEE Transactions on*, vol. 55, no. 12, pp. 5758–5772, 2009.
- [34] K. Rad, “Nearly sharp sufficient conditions on exact sparsity pattern recovery,” *Information Theory, IEEE Transactions on*, vol. 57, no. 7, pp. 4672–4679, 2011.
- [35] G. Reeves and M. Gastpar, “The sampling rate-distortion tradeoff for sparsity pattern recovery in compressed sensing,” *Information Theory, IEEE Transactions on*, vol. 58, no. 5, pp. 3065–3092, 2012.
- [36] —, “Approximate sparsity pattern recovery: Information-theoretic lower bounds,” *Information Theory, IEEE Transactions on*, vol. 59, no. 6, pp. 3451–3465, 2013.
- [37] G. Tang and A. Nehorai, “Performance analysis for sparse support recovery,” *Information Theory, IEEE Transactions on*, vol. 56, no. 3, pp. 1383–1399, 2010.
- [38] A. Tulino, G. Caire, S. Verdu, and S. Shamai, “Support recovery with sparsely sampled free random matrices,” *Information Theory, IEEE Transactions on*, vol. 59, no. 7, pp. 4243–4271, 2013.
- [39] W. Wang, M. Wainwright, and K. Ramchandran, “Information-theoretic limits on sparse signal recovery: Dense versus sparse measurement matrices,” *Information Theory, IEEE Transactions on*, vol. 56, no. 6, pp. 2967–2979, 2010.
- [40] M. Wainwright, “Sharp thresholds for high-dimensional and noisy sparsity recovery using  $l_1$ -constrained quadratic programming (lasso),” *Information Theory, IEEE Transactions on*, vol. 55, no. 5, pp. 2183–2202, 2009.
- [41] —, “Information-theoretic limits on sparsity recovery in the high-dimensional and noisy setting,” *Information Theory, IEEE Transactions on*, vol. 55, no. 12, pp. 5728–5741, 2009.
- [42] P. Zhao and B. Yu, “On model selection consistency of lasso,” *Machine Learning Research, Journal of*, vol. 7, pp. 2541–2563, December 2006.
- [43] C. Aksoylar, G. Atia, and V. Saligrama, “Sparse signal processing with linear and non-linear observations: A unified shannon theoretic approach,” *submitted to IEEE Trans. Info. Theory*, 2013. [Online]. Available: <http://arxiv.org/abs/1304.0682>
- [44] R. Calderbank and S. Jafarpour, “Finding needles in compressed haystacks,” in *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*. IEEE, 2012, pp. 3441–3444.
- [45] A. S. Bandeira, D. G. Mixon, and B. Recht, “Compressive classification and the rare eclipse problem,” *ArXiv e-prints*, Apr. 2014.
- [46] H. Reboredo, F. Renna, R. Calderbank, and M. R. Rodrigues, “Compressive classification of a mixture of gaussians: Analysis, designs and geometrical interpretation,” *submitted to IEEE Trans. Info. Theory*, 2014.
- [47] F. Renna, R. Calderbank, L. Carin, and M. R. Rodrigues, “Reconstruction of signals drawn from a gaussian mixture from noisy compressive measurements: Mmse phase transitions and beyond,” *submitted to IEEE Trans. Info. Theory*, 2013.
- [48] N. Santhanam and M. J. Wainwright, “Information-theoretic limits of graphical model selection in high dimensions,” in *Proc. Int. Symp. Information Theory (ISIT)*. IEEE, 2008, pp. 2136–2140.
- [49] J. Acharya, H. Das, and A. Orlitsky, “Tight bounds on profile redundancy and distinguishability,” in *Advances in Neural Information Processing Systems*, 2012, pp. 3266–3274.
- [50] E. Tuncel and D. Gunduz, “Identification and lossy reconstruction in noisy databases,” in *Information Theory Proceedings (ISIT), 2010 IEEE International Symposium on*. IEEE, 2010, pp. 191–195.
- [51] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. New York, NY: John Wiley and Sons, Ltd., 2006.
- [52] S. Verdu and T. Han, “A general formula for channel capacity,” *IEEE Trans. Info. Theory*, vol. 40, no. 4, pp. 1147–1157, 1994.
- [53] A. Bhattacharyya, “On a measure of divergence between two multinomial populations,” *The Indian Journal of Statistics*, vol. 7, no. 4, pp. 401–406, 1946.
- [54] T. Kailath, “The divergence and bhattacharyya distance measures in signal selection,” *IEEE Trans. Communication Technology*, vol. 15, no. 1, pp. 52–60, February 1967.
- [55] J. W. Silverstein, “The smallest eigenvalue of a large dimensional Wishart matrix,” *The Annals of Probability*, vol. 13, no. 4, pp. 1364–3468, 1985.
- [56] R. A. Horn and C. R. Johnson, *Matrix Analysis*, 2nd ed. Cambridge University Press, 2012.
- [57] R. Basri and D. W. Jacobs, “Lambertian reflectance and linear subspaces,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, no. 2, pp. 218–233, 2003.