

# Supplementary Information: Supervised Dimensionality Reduction for Big Data

Joshua T. Vogelstein<sup>1\*†</sup>, Eric W. Bridgeford<sup>1\*</sup>, Minh Tang<sup>1</sup>, Da Zheng<sup>1</sup>, Christopher Douville<sup>1</sup>, Randal Burns<sup>1</sup>, Mauro Maggioni<sup>1</sup>

<sup>1</sup> Johns Hopkins University 3900 N Charles Street, Baltimore, MD 21218, \*These authors contributed equally, <sup>†</sup>Corresponding Author, jovo@jhu.edu

## Supplementary Note 1 Theoretical Background

### 1.1 The Classification Problem

Let  $(\mathbf{X}, Y)$  be a pair of random variables, jointly sampled from  $F := F_{\mathbf{X}, Y} = F_{\mathbf{X}|Y}F_Y$ , with density denoted  $f_{\mathbf{X}, Y}$ . Let  $\mathbf{X}$  be a multivariate vector-valued random variable, such that its realizations live in  $p$  dimensional Euclidean space,  $\mathbf{x} \in \mathbb{R}^p$ . Let  $Y$  be a categorical random variable, whose realizations are discrete,  $y \in \{0, \dots, C-1\}$ . The goal of a classification problem is to find a function  $g(\mathbf{x})$  such that its output tends to be the true class label  $y$ :

$$g^*(\mathbf{x}) := \operatorname{argmax}_{g \in \mathcal{G}} \mathbb{P}[g(\mathbf{x}) = y]. \quad (1)$$

When the joint distribution of the data is known, then the Bayes optimal solution is:

$$g^*(\mathbf{x}) := \operatorname{argmax}_y f_{y, \mathbf{x}} = \operatorname{argmax}_y f_{\mathbf{x}|y} f_y = \operatorname{argmax}_y \{\log f_{\mathbf{x}|y} + \log f_y\} \quad (2)$$

Denote expected misclassification rate of classifier  $g$  for a given joint distribution  $F$ ,

$$L_g^F := \mathbb{E}[g(\mathbf{x}) \neq y] := \int \mathbb{P}[g(\mathbf{x}) \neq y] f_{\mathbf{x}, y} d\mathbf{x} dy, \quad (3)$$

where  $\mathbb{E}$  is the expectation, which in this case, is with respect to  $F_{\mathbf{X}, Y}$ . For brevity, we often simply write  $L_g$ , and we define  $L_* := L_{g^*}$ .

### 1.2 Linear Discriminant Analysis (LDA)

Linear Discriminant Analysis (LDA) is an approach to classification that uses a linear function of the first two moments of the distribution of the data. More specifically, let  $\boldsymbol{\mu}_j = \mathbb{E}[F_{\mathbf{X}|Y=j}]$  denote the class conditional mean, and let  $\boldsymbol{\Sigma} = \mathbb{E}[F_{\mathbf{X}}^2]$  denote the joint covariance matrix, and the class priors are  $\pi_j = \mathbb{P}[Y = j]$ . Using this notation, we can define the LDA classifier:

$$g_{\text{LDA}}(\mathbf{x}) := \operatorname{argmin}_y \left[ \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_y)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) - \log \pi_y \right], \quad (4)$$

Let  $L_{\text{LDA}}^F$  be the expected misclassification rate of the above classifier for distribution  $F$ . Assuming equal class prior and centered means, re-arranging a bit, we obtain

$$g_{\text{LDA}}(\mathbf{x}) := \operatorname{argmin}_y \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_y. \quad (5)$$

In words, the LDA classifier chooses the class that maximizes the magnitude of the projection of an input vector  $\mathbf{x}$  onto  $\Sigma^{-1}\boldsymbol{\mu}_y$ . When there are only two classes, letting  $\boldsymbol{\delta} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ , the above further simplifies to

$$g_{2\text{LDA}}(\mathbf{x}) := \mathbb{I}\{\mathbf{x}^\top \Sigma^{-1} \boldsymbol{\delta} > 0\}. \quad (6)$$

Note that the equal class prior and centered means assumptions merely changes the threshold constant from 0 to some other constant.

### 1.3 LDA Model

A statistical model is a family of distributions indexed by a parameter  $\boldsymbol{\theta} \in \Theta$ ,  $\mathcal{F}_\theta = \{F_\theta : \boldsymbol{\theta} \in \Theta\}$ . Consider the special case of the above where  $F_{X|Y=y}$  is a multivariate Gaussian distribution,  $\mathcal{N}(\boldsymbol{\mu}_y, \Sigma)$ , where each class has its own mean, but all classes have the same covariance. We refer to this model as the LDA model. Let  $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\mu}, \Sigma)$ , and let  $\Theta_{C\text{-LDA}} = (\Delta_C, \mathbb{R}^{p \times C}, \mathbb{R}_{>0}^{p \times p})$ , where  $\boldsymbol{\mu} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_C)$ ,  $\Delta_C$  is the  $C$  dimensional simplex, that is  $\Delta_C = \{\mathbf{x} : x_i \geq 0 \forall i, \sum_i x_i = 1\}$ , and  $\mathbb{R}_{>0}^{p \times p}$  is the set of positive definite  $p \times p$  matrices. Denote  $\mathcal{F}_{\text{LDA}} = \{F_\theta : \boldsymbol{\theta} \in \Theta_{\text{LDA}}\}$ . The following lemma is well known [5]:

**Lemma 1.**  $L_{\text{LDA}}^F = L_*^F$  for any  $F \in \mathcal{F}_{\text{LDA}}$ .

### Supplementary Note 2 Formal Definition of LOL and Related Classifiers

Let  $\mathbf{A} \in \mathbb{R}^{d \times p}$  be a ‘‘projection matrix’’, that is, a matrix that projects  $p$ -dimensional data into a  $d$ -dimensional subspace. The question that motivated this work is: what is the best projection matrix that we can estimate, to use to ‘‘pre-process’’ the data prior to classifying the data? Projecting the data  $\mathbf{x}$  onto a low-dimensional subspace, and then classifying via LDA in that subspace is equivalent to redefining the parameters in the low-dimensional subspace,  $\Sigma_A = \mathbf{A}\Sigma\mathbf{A}^\top \in \mathbb{R}^{d \times d}$  and  $\boldsymbol{\delta}_A = \mathbf{A}\boldsymbol{\delta} \in \mathbb{R}^d$ , and then using  $g_{\text{LDA}}$  in the low-dimensional space. When  $C = 2$ ,  $\pi_0 = \pi_1$ , and  $(\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1)/2 = \mathbf{0}$ , this amounts to:

$$g_A^d(\mathbf{x}) := \mathbb{I}\{(\mathbf{A}\mathbf{x})^\top \Sigma_A^{-1} \boldsymbol{\delta}_A > 0\}, \text{ where } \mathbf{A} \in \mathbb{R}^{d \times p}. \quad (7)$$

Let  $L_A^d := \int \mathbb{P}[g_A^d(\mathbf{x}) = y] f_{\mathbf{x},y} d\mathbf{x}dy$ . Our goal therefore is to be able to choose  $A$  for a given parameter setting  $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\delta}, \Sigma)$ , such that  $L_A$  is as small as possible (note that  $L_A$  will never be smaller than  $L_*$ ).

In the naive case where  $(\Sigma_A, \boldsymbol{\delta}_A)$  are known, we seek to solve the following linear optimization problem:

$$\begin{aligned} & \underset{\mathbf{A}}{\text{minimize}} && \mathbb{E}[\mathbb{I}\{\mathbf{x}^\top \mathbf{A}^\top \Sigma_A^{-1} \boldsymbol{\delta}_A > 0\} \neq y] \\ & \text{subject to} && \mathbf{A} \in \mathbb{R}^{d \times p}. \end{aligned} \quad (8)$$

When  $(\Sigma_A, \boldsymbol{\delta}_A)$  are not known, however, the optimization problem becomes non-convex. With  $\Sigma_A$  and  $\boldsymbol{\delta}_A$  as above:

$$\begin{aligned} & \underset{\mathbf{A}, \Sigma, \boldsymbol{\delta}}{\text{minimize}} && \mathbb{E}[\mathbb{I}\{\mathbf{x}^\top \mathbf{A}^\top \Sigma_A^{-1} \boldsymbol{\delta}_A > 0\} \neq y] \\ & \text{subject to} && \mathbf{A} \in \mathbb{R}^{d \times p}. \end{aligned} \quad (9)$$

While there are numerous approaches to solve related convex optimization problems through various sets of assumptions [10, 15], we do not consider such techniques in this manuscript theoretically. This

is because assuming either a structure for  $\Sigma_A$  or  $\delta_A$  presupposes an understanding of the properties of the feature space for wide data, which is often unsuitable if the dataset is large or has considerable complexity.

Let  $\mathcal{A} = \{\mathbf{A} : \mathbf{A} \in \mathbb{R}^{k \times p}, k \leq p\}$ ,  $\mathcal{A}^d = \{\mathbf{A} : \mathbf{A} \in \mathbb{R}^{d \times p}\}$ , and let  $\mathcal{A}_* \subset \mathcal{A}$  be the set of  $\mathbf{A}$  that minimizes Supp. Eq. (9), and let  $\mathbf{A}_* \in \mathcal{A}_*$ . Let  $L_{\mathbf{A}_*}^* = L_{\mathcal{A}_*}^*$  be the misclassification rate for any  $\mathbf{A} \in \mathcal{A}_*$ , that is,  $L_{\mathbf{A}_*}^*$  is the Bayes optimal misclassification rate for the classifier that composes  $\mathbf{A}$  with LDA.

In our opinion, Supp. Eq. (9) is the simplest supervised manifold learning problem there is: a two-class classification problem, where the data are multivariate Gaussians with shared, unknown covariances, the manifold is linear, and the classification is done via LDA. Nonetheless, solving Supp. Eq. (9) is difficult, because we do not know how to evaluate the integral analytically, and we do not know any algorithms that are guaranteed to find the global optimum in finite time. We proceed by studying a few natural choices for  $\mathbf{A}$ .

## 2.1 Bayes Optimal Projection

**Lemma 2.**  $\delta^\top \Sigma^{-1} \in \mathcal{A}_*$

*Proof.* Let  $\mathbf{B} = (\Sigma^{-1} \delta)^\top = \delta^\top (\Sigma^{-1})^\top = \delta^\top \Sigma^{-1}$ , so that  $\mathbf{B}^\top = \Sigma^{-1} \delta$ , and plugging this in to Supp. Eq. (7). By the above, and noting the symmetry and invertibility of  $\Sigma$ :

$$\begin{aligned} \Sigma_B &= \mathbf{B} \Sigma \mathbf{B}^\top = \delta^\top \Sigma^{-1} \Sigma (\delta^\top \Sigma^{-1})^\top \\ &= \delta^\top \Sigma^{-1} \Sigma \Sigma^{-1} \delta = \delta^\top \Sigma^{-1} \delta \\ \Rightarrow \Sigma_B^{-1} &= \delta^{-1} \Sigma \delta^{\top -1} \\ \delta_B &= \mathbf{B} \delta = \delta^\top \Sigma^{-1} \delta \end{aligned}$$

We obtain:

$$\begin{aligned} g_B(x) &= \mathbb{I}\{x^\top \mathbf{B}^\top \Sigma_B^{-1} \delta_B > 0\} \\ &= \mathbb{I}\{x^\top (\Sigma^{-1} \delta) (\Sigma_B^{-1} \delta_B) > 0\} && \text{plugging in } \mathbf{B} \\ &= \mathbb{I}\{x^\top (\Sigma^{-1} \delta) (\delta^{-1} \Sigma \delta^{\top -1} \delta^\top \Sigma^{-1} \delta) > 0\} && \text{plug in } \Sigma_B, \delta_B \text{ from above} \\ &= \mathbb{I}\{x^\top \Sigma^{-1} \delta > 0\} \end{aligned}$$

In other words, letting  $\mathbf{B}$  be the Bayes optimal projection recovers the Bayes classifier, as it should. Or, more formally, for any  $F \in \mathcal{F}_{\text{LDA}}$ ,  $L_{\delta^\top \Sigma^{-1}} = L_*$ .  $\square$

## 2.2 Principle Components Analysis (PCA) Projection

Principle Components Analysis (PCA) finds the directions of maximal variance in a dataset. PCA is closely related to eigendecompositions and singular value decompositions (SVD). In particular, the top left singular vector of a matrix  $\mathbf{X} \in \mathbb{R}^{p \times n}$ , whose columns are centered, is the eigenvector with the largest eigenvalue of the centered covariance matrix  $\mathbf{X} \mathbf{X}^\top$ . SVD enables one to estimate this eigenvector without ever forming the outer product matrix, because SVD factorizes a matrix  $\mathbf{X}$  into  $\mathbf{U} \mathbf{S} \mathbf{V}^\top$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal  $p \times p$  matrices, and  $\mathbf{S}$  is a diagonal matrix, whose diagonal values are decreasing,  $s_1 \geq s_2 \geq \dots > s_n$ . Defining  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p]$ , where each  $\mathbf{u}_i \in \mathbb{R}^p$ , then  $\mathbf{u}_i$  is the  $i^{\text{th}}$  eigenvector, and  $s_i$  is the square root of the  $i^{\text{th}}$  eigenvalue of  $\mathbf{X} \mathbf{X}^\top$ . Let  $\mathbf{A}_d^{\text{PCA}} =$

$[\mathbf{u}_1, \dots, \mathbf{u}_d]$  be the truncated PCA orthonormal matrix, and let  $\mathbf{I}_{d \times p}$  denote a  $d \times p$  dimensional identity matrix.

The PCA matrix is perhaps the most obvious choice of an orthonormal matrix for several reasons. First, truncated PCA minimizes the squared error loss between the original data matrix and all possible rank  $d$  representations:

$$\operatorname{argmin}_{\mathbf{A} \in \mathbb{R}^{d \times p}} \|\mathbf{X} - \mathbf{A}^T \mathbf{A}\|_F^2. \quad (10)$$

Second, the ubiquity of PCA has led to a large number of highly optimized numerical libraries for computing PCA (for example, LAPACK [2]).

In this supervised setting, we consider two different variants of PCA, each based on centering the data differently. For the first one, which we refer to as ‘‘pooled PCA’’ (or just PCA for brevity), we center the data by subtracting the ‘‘pooled mean’’ from each sample, that is, we let  $\tilde{\mathbf{x}}_i = \mathbf{x} - \boldsymbol{\mu}$ , where  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}]$ . For the second, which we refer to as ‘‘class conditional PCA’’, we center the data by subtracting the ‘‘class-conditional mean’’ from each sample, that is, we let  $\tilde{\mathbf{x}}_i = \mathbf{x} - \boldsymbol{\mu}_y$ , where  $\boldsymbol{\mu}_y = \mathbb{E}[\mathbf{x}|Y = y]$ .

Notationally, let  $\mathbf{U}_d = [\mathbf{u}_1, \dots, \mathbf{u}_d] \in \mathbb{R}^{p \times d}$ , and note that  $\mathbf{U}_d^T \mathbf{U}_d = \mathbf{I}_{d \times d}$  and  $\mathbf{U}_d \mathbf{U}_d^T = \mathbf{I}_{p \times p}$ . Similarly, let  $\mathbf{U} \mathbf{S} \mathbf{U}^T = \boldsymbol{\Sigma}$ , and  $\mathbf{U} \mathbf{S}^{-1} \mathbf{U}^T = \boldsymbol{\Sigma}^{-1}$ . Let  $\mathbf{S}_d$  be the matrix whose diagonal entries are the eigenvalues, up to the  $d^{\text{th}}$  one, that is  $\mathbf{S}_d(i, j) = s_i$  for  $i = j \leq d$  and zero otherwise. Similarly,  $\boldsymbol{\Sigma}_d = \mathbf{U} \mathbf{S}_d \mathbf{U}^T = \mathbf{U}_d \mathbf{S}_d \mathbf{U}_d^T$ . Reduced-rank LDA (rrLDA) is a regularized LDA algorithm. Specifically, rather than using the full rank covariance matrix, it uses a rank- $d$  approximation. Formally, let  $g_{\text{LDA}} := \mathbb{I}\{\mathbf{x} \boldsymbol{\Sigma}^{-1} \boldsymbol{\delta} > 0\}$  be the LDA classifier, and let  $g_{\text{LDA}}^d := \mathbb{I}\{\mathbf{x} \boldsymbol{\Sigma}_d^{-1} \boldsymbol{\delta} > 0\}$  be the regularized LDA classifier, that is, the LDA classifier where the the bottom  $p - d$  eigenvalues of the covariance matrix are set to zero.

**Lemma 3.** *Using class-conditional PCA to pre-process the data, then using LDA on the projected data, is equivalent to rrLDA.*

*Proof.* Plugging  $\mathbf{U}_d$  into Eq. (7) for  $\mathbf{A}$ , and considering only the left side of the operand, we have

$$\begin{aligned} (\mathbf{A}\mathbf{x})^T \boldsymbol{\Sigma}_A^{-1} \boldsymbol{\delta}_A &= \mathbf{x}^T \mathbf{A}^T \mathbf{A} \boldsymbol{\Sigma}^{-1} \mathbf{A}^T \mathbf{A} \boldsymbol{\delta}, \\ &= \mathbf{x}^T \mathbf{U}_d \mathbf{U}_d^T \boldsymbol{\Sigma}^{-1} \mathbf{U}_d \mathbf{U}_d^T \boldsymbol{\delta}, \\ &= \mathbf{x}^T \mathbf{U}_d \mathbf{U}_d^T \mathbf{U} \mathbf{S}^{-1} \mathbf{U}^T \mathbf{U}_d \mathbf{U}_d^T \boldsymbol{\delta}, \\ &= \mathbf{x}^T \mathbf{U}_d \mathbf{I}_{d \times p} \mathbf{S}^{-1} \mathbf{I}_{p \times d} \mathbf{U}_d^T \boldsymbol{\delta}, \\ &= \mathbf{x}^T \mathbf{U}_d \mathbf{S}_d^{-1} \mathbf{U}_d^T \boldsymbol{\delta}, \\ &= \mathbf{x}^T \boldsymbol{\Sigma}_d^{-1} \boldsymbol{\delta}. \end{aligned}$$

□

The implication of this lemma is that if one desires to implement rrLDA, rather than first learning the eigenvectors and then learning LDA, one can instead directly implement regularized LDA by setting the bottom  $p - d$  eigenvalues to zero. This latter approach removes the requirement to run SVD twice, therefore reducing the computational burden as well as the possibility of numerical instability issues. We therefore refer to the projection composed of  $d$  eigenvectors of the class-conditionally centered covariance matrices,  $\mathbf{A}_{\text{LDA}}^d$ .

### 2.3 Linear Optimal Low-Rank (LOL) Projection

The basic idea of LOL is to use both  $\delta$  and the top  $d$  eigenvectors of the class-conditionally centered covariance. When there are only two classes,  $\delta = \mu_0 - \mu_1$ . When there are  $C > 2$  classes, there are  $\binom{C}{2} = \frac{C!}{2(C-2)!}$  pairwise combinations,  $\delta_{ij} = \mu_i - \mu_j$  for all  $i \neq j$ . However, since  $\binom{C}{2}$  is nearly  $C^2$ , when  $C$  is large, this would mean incorporating many mean difference vectors. Note that  $[\delta_{1,2}, \delta_{1,3}, \dots, \delta_{C-1,C}]$  is in fact a rank  $C - 1$  matrix, because it is a linear function of the  $C$  different means. Therefore, we only need  $C - 1$  differences to span the space of all pairwise differences. To mitigate numerical instability issues, we adopt the following convention. For each class, estimate the expected mean and the number of samples per class,  $\mu_c$  and  $\pi_c$ . Sort the means in order of decreasing  $\pi_c$ , so that  $\pi_{(1)} > \pi_{(2)} > \dots > \pi_{(C)}$ . Then, subtract  $\mu_{(1)}$  from all other means:  $\delta_i = \mu_{(1)} - \mu_{(i)}$ , for  $i = 2, \dots, C$ . Finally,  $\delta = [\delta_2, \dots, \delta_C]$ .

Given  $\delta$  and  $\mathbf{A}_{\text{LDA}}^{d-1}$ , to obtain LOL naïvely, we could simply concatenate the two,  $\mathbf{A}_{\text{LOL},naive}^d = [\delta, \mathbf{A}_{\text{LDA}}^{d-1}]$ . Recall that eigenvectors are orthonormal. To maintain orthonormality between the eigenvectors and vectors of  $\delta$ , we could easily apply Gram-Schmidt,  $\mathbf{A}_{\text{LOL},naive}^d = \text{Orth}([\delta, \mathbf{A}_{\text{LDA}}^{d-1}])$ . In practice, this orthogonalization step does not matter much, so we ignore it hereafter. To ensure that  $\delta$  and  $\Sigma$  are balanced appropriately, we normalize each vector in  $\delta$  to have norm unity. Formally, let  $\tilde{\delta}_j = \delta_j / \|\delta_j\|$ , where  $\delta_j$  is the  $j^{\text{th}}$  difference of the mean vector and let  $\mathbf{A}_{\text{LOL}}^d = [\tilde{\delta}, \mathbf{A}_{\text{LDA}}^{d-(C-1)}]$ .

When the distribution of the data is not provided, each of the above terms must be estimated from the data. We use the maximum likelihood estimators for each, specifically:

$$n_c = \sum_{i=1}^n \mathbb{I}\{y_i = c\}, \quad (11)$$

$$\hat{\pi}_c = \frac{n_c}{n}, \quad (12)$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad (13)$$

$$\hat{\mu}_c = \frac{1}{n_c} \sum_{i=1}^n \mathbf{x}_i \mathbb{I}\{y_i = c\}. \quad (14)$$

For completeness, below we provide pseudocode for learning the sample version of LOL. The population version does not require the estimation of the parameters.

### 2.4 rrLDA is rotationally invariant

For certain classification tasks, the observed dimensions (or features) have intrinsic value, e.g. when simple interpretability is desired. However, in many other contexts, interpretability is less important [3]. When the exploitation task at hand is invariant to rotations, then we have no reason to restrict our search space to be sparse in the observed dimensions. For example, we can consider sparsity in the eigenvector basis. Let  $\mathbf{W}$  be a rotation matrix, that is  $\mathbf{W} \in \mathcal{W} = \{\mathbf{W} : \mathbf{W}^T = \mathbf{W}^{-1} \text{ and } \det(\mathbf{W}) = 1\}$ . Moreover, let  $\mathbf{W} \circ F$  denote the distribution  $F$  after transformation by an operator  $\mathbf{W}$ . For example, if  $F = \mathcal{N}(\mu, \Sigma)$  then  $\mathbf{W} \circ F = \mathcal{N}(\mathbf{W}\mu, \mathbf{W}\Sigma\mathbf{W}^T)$ .

**Definition 1.** A rotationally invariant classifier has the following property:

$$L_g^F = L_g^{\mathbf{W} \circ F}, \quad F \in \mathcal{F} \text{ and } \mathbf{W} \in \mathcal{W}. \quad (15)$$

In words, the Bayes risk of using classifier  $g$  on distribution  $F$  is unchanged if  $F$  is first rotated.

Now, we can state the main lemma of this subsection:  $\text{rrLDA}$  is rotationally invariant.

**Lemma 4.**  $L_{\text{LDA}}^F = L_{\text{LDA}}^{W \circ F}$ , for any  $F \in \mathcal{F}$ .

*Proof.*  $\text{rrLDA}$  is in fact simply thresholding  $x^\top \Sigma^{-1} \delta$  whenever its value is larger than some constant. Thus, we can demonstrate rotational invariance by demonstrating that  $x^\top \Sigma^{-1} \delta$  is rotationally invariant.

$$\begin{aligned}
(Wx)^\top (W\Sigma W^\top)^{-1} W\delta &= x^\top W^\top (WUSU^\top W^\top)^{-1} W\delta && \text{by substituting } USU^\top \text{ for } \Sigma \\
&= x^\top W^\top (\tilde{U}S\tilde{U}^\top)^{-1} W\delta && \text{by letting } \tilde{U} = WU \\
&= x^\top W^\top (\tilde{U}S^{-1}\tilde{U}^\top)W\delta && \text{by the laws of matrix inverse} \\
&= x^\top W^\top WUS^{-1}U^\top W^\top W\delta && \text{by un-substituting } WU = \tilde{U} \\
&= x^\top US^{-1}U^\top \delta && \text{because } W^\top W = I \\
&= x^\top \Sigma^{-1} \delta && \text{by un-substituting } US^{-1}U^\top = \Sigma
\end{aligned}$$

□

One implication of this lemma is that we can reparameterize without loss of generality. Specifically, defining  $W := U^\top$  yields a change of variables:  $\Sigma \mapsto S$  and  $\delta \mapsto U^\top \delta := \delta'$ , where  $S$  is a diagonal covariance matrix. Moreover, let  $d = (\sigma_1, \dots, \sigma_D)^\top$  be the vector of eigenvalues, then  $S^{-1} \delta' = d^{-1} \odot \tilde{\delta}$ , where  $\odot$  is the Hadamard (entrywise) product. The  $\text{LDA}$  classifier may therefore be encoded by a unit vector,  $\tilde{d} := \frac{1}{m} d^{-1} \odot \tilde{\delta}'$ , and its magnitude,  $m := \left\| d^{-1} \odot \tilde{\delta}' \right\|$ . This will be useful later.

## 2.5 Rotation of Projection Based Linear Classifiers

By a similar argument as above, one can easily show that:

$$(AWx)^\top (AW\Sigma W^\top A^\top)^{-1} AW\delta = x^\top (W^\top A^\top)(AW)\Sigma^{-1}(W^\top A^\top)(AW)\delta \quad (16)$$

$$= x^\top Y^\top Y\Sigma^{-1}Y^\top Y\delta \quad (17)$$

$$= x^\top Z\Sigma^{-1}Z^\top \delta \quad (18)$$

$$= x^\top (Z\Sigma Z^\top)^{-1} \delta = x^\top \tilde{\Sigma}_d^{-1} \delta, \quad (19)$$

where  $Y = AW \in \mathbb{R}^{d \times p}$  so that  $Z = Y^\top Y$  is a symmetric  $p \times p$  matrix of rank  $d$ . In other words, rotating and then projecting is equivalent to a change of basis. The implications of the above is:

**Lemma 5.**  $g_A$  is rotationally invariant if and only if  $\text{span}(A) = \text{span}(\Sigma_d)$ . In other words,  $\text{rrLDA}$  is the only rotationally invariant projection.

## 2.6 Low-Rank Canonical Correlation Analysis

We now contrast `LOL` and low-rank CCA. For discriminant analysis, low-rank CCA corresponds to finding the eigenvectors of  $S_X^\dagger S_B$  where

$$S_X = \sum_i (X_i - \bar{X})(X_i - \bar{X})^\top; \quad \bar{X} = \sum_i X_i \quad (20)$$

is the sample covariance matrix of the  $X_i$ ,  $S_X^\dagger$  is the inverse of  $S_X$  (or Moore-Penrose pseudo-inverse of  $S_X$  if  $S_X$  is not invertible), and

$$S_B = \frac{n_0}{n}(\bar{X}_0 - \bar{X})(\bar{X}_0 - \bar{X})^\top + \frac{n_1}{n}(\bar{X}_1 - \bar{X})(\bar{X}_1 - \bar{X})^\top; \quad \bar{X}_j = \sum_{i: Y_i=j} X_i \text{ for } j \in \{0, 1\} \quad (21)$$

is the between class covariance matrix [18]. It is widely known (see section 11.5 of [16]) that if  $S_X$  is invertible then the above formulation reduces to that of Fisher L, namely that of finding  $\hat{v}$  satisfying

$$\hat{v} = \operatorname{argmax}_{v \neq 0} \frac{v^\top S_B v}{v^\top S_W v} \quad (22)$$

$$S_W = \sum_{i: Y_i=0} (X_i - \bar{X}_0)(X_i - \bar{X}_0)^\top + \sum_{i: Y_i=1} (X_i - \bar{X}_1)(X_i - \bar{X}_1)^\top; \quad (23)$$

where  $S_W$  is the pooled within-sample covariance matrix and  $S_X = S_W + S_B$ . In the context of our current paper where  $X$  is assumed to be high-dimensional, it is well-known that  $S_X$  is not a good estimator of the population covariance matrix  $\Sigma_X = \mathbb{E}[(X - \mu)(X - \mu)^\top]$  and thus computing  $S_X^{-1}$  is suboptimal for subsequent inference unless some form of regularization is employed. Our consideration of low-rank linear transformations  $AX$  provides one principled approach to regularizations of high-dimensional  $S_X$ . In contrast, the above (unregularized) formulation of low-rank CCA frequently yields discrimination direction vectors corresponding to “maximum data piling” (MDP) directions [1, 18] in high-dimensional settings (and always yield maximum data piling directions when  $p \geq n$ ). These MDP directions lead to *perfect* discrimination of the training data, but can suffer from poor generalization performance, as the examples in [1, 18] indicate.

Naïvely computing the low-rank CCA projection requires storing and inverting a  $p \times p$  matrix. However, we devised an implementation for low-rank CCA that does not require ever materializing this matrix. Modern eigensolvers compute eigenvalues by performing a sequence of matrix vector multiplication. For example, to compute eigenvalues of  $S_X$ , an eigensolver performs  $S_X v$  multiple times until the algorithm converges. Assume that the number of iteration is  $i$ , the computation complexity of the eigensolver is  $O(n \times p \times i)$ . Performing pseudo-inverse of  $S_X$  computes truncated SVD on  $S_X$ , resulting in  $S_X v = \sum_i (X_i - \bar{X})((X_i - \bar{X})^\top v)$ . Here we never physically generate  $S_X$ . Instead, we always compute  $v' = (X_i - \bar{X})^\top v$  and then  $v'' = (X_i - \bar{X})v'$  to compute  $S_X v$ . Assume  $k$  classes,  $S_X v$  has the computation complexity of  $O(n \times p \times k)$  and the space complexity of  $O(n \times p \times k)$ .  $S_X$  can be decomposed into  $U\Sigma V$ , where  $U$  is a  $n \times n$  matrix and  $V$  is a  $n \times p$  matrix.

$$S_X^\dagger S_B = U\Sigma^{-1}V\left(\frac{n_0}{n}(\bar{X}_0 - \bar{X})(\bar{X}_0 - \bar{X})^\top + \frac{n_1}{n}(\bar{X}_1 - \bar{X})(\bar{X}_1 - \bar{X})^\top\right). \quad (24)$$

Computing eigenvalues of  $S_X^\dagger S_B$  requires

$$S_X^\dagger S_B v = U\Sigma^{-1}V\left(\frac{n_0}{n}(\bar{X}_0 - \bar{X})((\bar{X}_0 - \bar{X})^\top v) + \frac{n_1}{n}(\bar{X}_1 - \bar{X})((\bar{X}_1 - \bar{X})^\top v)\right). \quad (25)$$

Similar to  $S_X v$ , we never physically generate  $S_X^\dagger$  or  $S_B$ . Instead, we always multiply the terms on the right with  $v$  first, which results in the computation complexity of  $O(n \times p)$  and the space complexity of  $O(n \times p)$ . To our knowledge, this algorithm is novel, and the implementation is also of course novel.

## Supplementary Note 3 Simulations

Let  $f_{x|y}$  denote the conditional distribution of  $X$  given  $Y$ , and let  $f_y$  denote the prior probability of  $Y$ . For simplicity, assume that realizations of the random variable  $X$  are  $p$ -dimensional vectors,  $x \in \mathbb{R}^p$ , and realizations of the random variable  $Y$  are binary,  $y \in \{0, 1\}$ . For most simulation settings, each class is Gaussian:  $f_{x|y} = \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$ , where  $\boldsymbol{\mu}_y$  is the class-conditional mean and  $\boldsymbol{\Sigma}_y$  is the class-conditional covariance. Moreover, we assume  $f_y$  is a Bernoulli distribution with probability  $\pi$  that  $y = 1$ ,  $f_y = \mathcal{B}(\pi)$ . We typically assume that both classes are equally likely,  $\pi = 0.5$ , and the covariance matrices are the same,  $\boldsymbol{\Sigma}_0 = \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}$ . Under such assumptions, we merely specify  $\boldsymbol{\theta} = \{\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}\}$ . We consider the following simulation settings:

### Stacked Cigars

- $\boldsymbol{\mu}_0 = \mathbf{0}$ ,
- $\boldsymbol{\mu}_1 = (a, b, a, \dots, a)$ ,
- $\boldsymbol{\Sigma}$  is a diagonal matrix, with diagonal vector,  $\mathbf{d} = (1, b, 1, \dots, 1)$ ,

where  $a = 0.15$  and  $b = 4$ .

### Trunk

- $\boldsymbol{\mu}_0 = b/\sqrt{(1, 3, 5, \dots, 2p)}$ ,
- $\boldsymbol{\mu}_1 = -\boldsymbol{\mu}_0$ ,
- $\boldsymbol{\Sigma}$  is a diagonal matrix, with diagonal vector,  $\mathbf{d} = 100/\sqrt{(p, p-1, p-2, \dots, 1)}$ ,

where  $b = 4$ .

**Rotated Trunk** Same as Trunk, but the data are randomly rotated, that is, we sample  $Q$  uniformly from the set of  $p$ -dimensional rotation matrices, and then set:

- $\boldsymbol{\mu}_0 \leftarrow Q\boldsymbol{\mu}_0$ ,
- $\boldsymbol{\mu}_1 \leftarrow Q\boldsymbol{\mu}_1$ ,
- $\boldsymbol{\Sigma} \leftarrow Q\boldsymbol{\Sigma}Q^\top$ .

**3 Classes** Same as Trunk, but with a third mean equal to the zero vector,  $\boldsymbol{\mu}_2 = \mathbf{0}$ .

- $\boldsymbol{\mu}_0 = b/\sqrt{(1, 3, 5, \dots, 2p)}$ ,
- $\boldsymbol{\mu}_1 = -\boldsymbol{\mu}_0$ ,
- $\boldsymbol{\mu}_2 = \mathbf{0}$ ,
- $\boldsymbol{\Sigma}$  is a diagonal matrix, with diagonal vector,  $\mathbf{d} = 100/\sqrt{(p, p-1, p-2, \dots, 1)}$ ,

where  $b = 4$ .

**Robust** An experiment in which outliers are present for estimation of the projection matrix, but removed for training and testing of the classifier. This is due to the strong amount of noise present in the robust experiment will lead to poor generalizability of the estimated LDA classifier. Parameters indexed by  $i$  correspond to the generative model for the inliers, and those with  $o$  correspond to the outliers.

- $\boldsymbol{\mu}_0^{(i)} = b/\sqrt{(1, 3, 5, \dots, p)}$  for the first  $p/2$  dimensions and 0 otherwise,



- $\boldsymbol{\mu}_1^{(i)} = -\boldsymbol{\mu}_0$ ,
- $\boldsymbol{\Sigma}^{(i)} = b^3 / \sqrt{(1, 2, \dots, p)}$ ,
- $\boldsymbol{\mu}^{(o)} = \mathbf{0}$ ,
- $\boldsymbol{\Sigma}^{(o)} = b^6 / \sqrt{(1, 2, \dots, p)}$ ,
- $\pi^{(i)} = 0.7$ ,
- $\pi^{(o)} = 0.3$ ,

and outliers are randomly assigned class 0 or class 1 with equal probability.

**Cross** An experiment in which the two classes have identical means but different covariance matrices, meaning the optimal discriminant boundary is quadratic.

- $\boldsymbol{\mu}_0 = \boldsymbol{\mu}_1 = \mathbf{0}$ ,
- $\Sigma_0$  is a diagonal matrix, with diagonal  $(a, \dots, a, b, \dots, b)$  where the first  $\frac{d}{3}$  elements are  $a$ , and the rest are  $b$ ,
- $\Sigma_1$  is a diagonal matrix, with diagonal  $(b, \dots, b, a, \dots, a, b, \dots, b)$  where the middle  $\frac{d}{3}$  elements are  $a$ , and the others are  $b$ ,

and we let  $a = 1$ , and  $b = \frac{1}{4}$ .

**Hump- $K$**  An experiment with  $K$  classes, in which the class means display an alternating series of humps, and the class covariance is a scalar multiple of the identity.

- $\pi_k = \frac{1}{K}$
- $x_{l,k} = \lfloor -\frac{K}{2} \rfloor$  the left endpoint of the hump
- $x_{r,k} = d - x_{l, K-k+1}$  the right endpoint of the hump
- $x_{m,k} = \lfloor \frac{x_{l,k} + x_{r,k}}{2} \rfloor$  the midpoint of the hump
- Let  $a_k, b_k, c_k$  be the unique coefficients such that  $c_k + b_k x + a_k x^2$  passes through  $x_{l,k}$  at  $y = 0$ , passes through  $x_{r,k}$  at  $y = b$ , and passes through  $x_{m,k}$  at  $y = 0$ .
- Let  $\alpha_k = \begin{cases} 1 & k \text{ is odd} \\ -1 & k \text{ is even} \end{cases}$
- for  $j = 1, \dots, d$ , let  $\mu_{k,j} = \begin{cases} 0 & j \notin [x_{l,k}, x_{r,k}] \\ c_k + b_k j + a_k j^2 & j \in [x_{l,k}, x_{r,k}] \end{cases}$
- $\boldsymbol{\Sigma}$  is a diagonal matrix, with diagonal vector  $(\sigma, \dots, \sigma)$ .

where  $b = 4$  and  $\sigma = \frac{100}{K}$ .

**Computational Efficiency Experiments** These experiments used the Trunk setting, increasing the observed dimensionality.

**Hypothesis Testing Experiments** We considered two related joint distributions here. The first joint (Diagonal) is described by:

- $\boldsymbol{\mu}_0 = \mathbf{0}$ ,
- $\tilde{\boldsymbol{\mu}}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\boldsymbol{\mu}_1 = \tilde{\boldsymbol{\mu}}_1 / \|\tilde{\boldsymbol{\mu}}_1\|$ ,
- $\boldsymbol{\Sigma}$  is the same Toeplitz matrix (where the top row is  $\rho^{(0,1,2,\dots,p-1)}$ ), and the matrix is rescaled to have a Frobenius norm of 50.

The second (Dense) is the same except that the eigenvectors are uniformly random sampled orthonormal matrices, rather than the identity matrix.

**Regression Experiments** In this experiment we used a distribution similar to the Toeplitz distribution as described above, but  $y$  was a linear function of  $x$ , that is,  $y = Ax$ , where  $x \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , where  $\Sigma$  is the above described Toeplitz matrix, and  $A$  is a diagonal matrix whose first two diagonal elements are non-zero, and the rest are zero.

## Supplementary Note 4 Theorems and Proofs of Main Result

### 4.1 Chernoff information

We now introduce the notion of the Chernoff information, which serves as our surrogate measure for the Bayes error of any classification procedure given the *projected* data. Our discussion of the Chernoff information is under the context of decision rules for hypothesis testing, nevertheless, as evidenced by the fact that the *maximum a posteriori* decision rule—equivalently the Bayes classifier—achieves the Chernoff information rate, this distinction between hypothesis testing and classification is mainly for ease of exposition.

Let  $F_0$  and  $F_1$  be two absolutely continuous multivariate distributions in  $\Omega \subset \mathbb{R}^d$  with density functions  $f_0$  and  $f_1$ , respectively. Suppose that  $X_1, X_2, \dots, X_m$  are independent and identically distributed random variables, with  $X_i$  distributed either  $F_0$  or  $F_1$ . We are interested in testing the simple null hypothesis  $\mathbb{H}_0: F = F_0$  against the simple alternative hypothesis  $\mathbb{H}_1: F = F_1$ . A test  $T$  is a sequence of mapping  $T_m: \Omega^m \mapsto \{0, 1\}$  such that given  $X_1 = x_1, X_2 = x_2, \dots, X_m = x_m$ , the test rejects  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  if  $T_m(x_1, x_2, \dots, x_m) = 1$ ; similarly, the test decides  $\mathbb{H}_1$  instead of  $\mathbb{H}_0$  if  $T_m(x_1, x_2, \dots, x_m) = 0$ . The Neyman-Pearson lemma states that, given  $X_1 = x_1, X_2 = x_2, \dots, X_m = x_m$  and a threshold  $\eta_m \in \mathbb{R}$ , the likelihood ratio test rejects  $\mathbb{H}_0$  in favor of  $\mathbb{H}_1$  whenever

$$\left( \sum_{i=1}^m \log f_0(x_i) - \sum_{i=1}^m \log f_1(x_i) \right) \leq \eta_m. \quad (26)$$

Moreover, the likelihood ratio test is the most powerful test at significance level  $\alpha_m = \alpha(\eta_m)$ , i.e., the likelihood ratio test minimizes the type II error  $\beta_m$  subject to the constraint that the type I error is at most  $\alpha_m$ .

Assume that  $\pi \in (0, 1)$  is a prior probability of  $\mathbb{H}_0$  being true. Then, for a given  $\alpha_m^* \in (0, 1)$ , let  $\beta_m^* = \beta_m^*(\alpha_m^*)$  be the type II error associated with the likelihood ratio test when the type I error is at most  $\alpha_m^*$ . The quantity  $\inf_{\alpha_m^* \in (0, 1)} \pi \alpha_m^* + (1 - \pi) \beta_m^*$  is then the Bayes risk in deciding between  $\mathbb{H}_0$  and  $\mathbb{H}_1$  given the  $m$  independent random variables  $X_1, X_2, \dots, X_m$ . A classical result of Chernoff [6] states that the Bayes risk is intrinsically linked to a quantity known as the *Chernoff information*. More specifically, let  $C(F_0, F_1)$  be the quantity

$$\begin{aligned} C(F_0, F_1) &= -\log \left[ \inf_{t \in (0, 1)} \int_{\mathbb{R}^d} f_0^t(\mathbf{x}) f_1^{1-t}(\mathbf{x}) d\mathbf{x} \right] \\ &= \sup_{t \in (0, 1)} \left[ -\log \int_{\mathbb{R}^d} f_0^t(\mathbf{x}) f_1^{1-t}(\mathbf{x}) d\mathbf{x} \right] \end{aligned} \quad (27)$$

Then we have

$$\lim_{m \rightarrow \infty} \frac{1}{m} \inf_{\alpha_m^* \in (0, 1)} \log(\pi \alpha_m^* + (1 - \pi) \beta_m^*) = -C(F_0, F_1). \quad (28)$$

Thus  $C(F_0, F_1)$  is the *exponential* rate at which the Bayes error  $\inf_{\alpha_m^* \in (0,1)} \pi \alpha_m^* + (1 - \pi) \beta_m^*$  decreases as  $m \rightarrow \infty$ ; we also note that the  $C(F_0, F_1)$  is independent of  $\pi$ . We also define, for a given  $t \in (0, 1)$  the Chernoff divergence  $C_t(F_0, F_1)$  between  $F_0$  and  $F_1$  by

$$C_t(F_0, F_1) = -\log \int_{\mathbb{R}^d} f_0^t(\mathbf{x}) f_1^{1-t}(\mathbf{x}) d\mathbf{x}. \quad (29)$$

The Chernoff divergence is an example of a  $f$ -divergence as defined in [8]. When  $t = 1/2$ ,  $C_t(F_0, F_1)$  is the Bhattacharyya distance between  $F_0$  and  $F_1$ .

The result of Supp. Eq. (28) can be extended to  $K + 1 \geq 2$  hypothesis, with the exponential rate being the minimum of the Chernoff information between any pair of hypothesis. More specifically, let  $F_0, F_1, \dots, F_K$  be distributions on  $\mathbb{R}^d$  and let  $X_1, X_2, \dots, X_m$  be independent and identically distributed random variables with distribution  $F \in \{F_0, F_1, \dots, F_K\}$ . Our inference task is in determining the distribution of the  $X_i$  among the  $K+1$  hypothesis  $\mathbb{H}_0: F = F_0, \dots, \mathbb{H}_K: F = F_K$ . Suppose also that hypothesis  $\mathbb{H}_k$  has a *priori* probability  $\pi_k$ . For any decision rule  $g$ , the risk of  $g$  is  $r(g) = \sum_k \pi_k \sum_{l \neq k} \alpha_{lk}(g)$  where  $\alpha_{lk}(g)$  is the probability of accepting hypothesis  $\mathbb{H}_l$  when hypothesis  $\mathbb{H}_k$  is true. Then we have [13]

$$\inf_g \lim_{m \rightarrow \infty} \frac{r(g)}{m} = -\min_{k \neq l} C(F_k, F_l), \quad (30)$$

where the infimum is over all decision rules  $g$ , i.e., for any  $g$ ,  $r(g)$  decreases to 0 as  $m \rightarrow \infty$  at a rate no faster than  $\exp(-m \min_{k \neq l} C(F_k, F_l))$ .

When the distributions  $F_0$  and  $F_1$  are multivariate normal, that is,  $F_0 = \mathcal{N}(\mu_0, \Sigma_0)$  and  $F_1 = \mathcal{N}(\mu_1, \Sigma_1)$ ; then, denoting by  $\Sigma_t = t\Sigma_0 + (1 - t)\Sigma_1$ , we have

$$C(F_0, F_1) = \sup_{t \in (0,1)} \left( \frac{t(1-t)}{2} (\mu_1 - \mu_0)^\top \Sigma_t^{-1} (\mu_1 - \mu_0) + \frac{1}{2} \log \frac{|\Sigma_t|}{|\Sigma_0|^t |\Sigma_1|^{1-t}} \right). \quad (31)$$

## 4.2 Projecting data and Chernoff information

We now discuss how the Chernoff information characterizes the effect a linear transformation  $A$  of the data has on classification accuracy. We start with the following simple result whose proof follows directly from Supp. Eq. (30).

**Lemma 6.** *Let  $F_0 = \mathcal{N}(\mu_0, \Sigma)$  and  $F_1 \sim \mathcal{N}(\mu_1, \Sigma)$  be two multivariate normals with equal covariance matrices. For any linear transformation  $A$ , let  $F_0^{(A)}$  and  $F_1^{(A)}$  denote the distribution of  $AX$  when  $X \sim F_0$  and  $X \sim F_1$ , respectively. We then have*

$$\begin{aligned} C(F_0^{(A)}, F_1^{(A)}) &= \frac{1}{8} (\mu_1 - \mu_0)^\top A^\top (A \Sigma A^\top)^{-1} A (\mu_1 - \mu_0) \\ &= \frac{1}{8} (\mu_1 - \mu_0)^\top \Sigma^{-1/2} \Sigma^{1/2} A^\top (A \Sigma A^\top)^{-1} A \Sigma^{1/2} \Sigma^{-1/2} (\mu_1 - \mu_0) \\ &= \frac{1}{8} \|P_{\Sigma^{1/2} A^\top} \Sigma^{-1/2} (\mu_1 - \mu_0)\|_F^2 \end{aligned} \quad (32)$$

where  $P_Z = Z(Z^\top Z)^{-1} Z^\top$  denotes the matrix corresponding to the orthogonal projection onto the columns of  $Z$ .

Thus for a classification problem where  $X|Y = 0$  and  $X|Y = 1$  are distributed multivariate normals with mean  $\mu_0$  and  $\mu_1$  and the same covariance matrix  $\Sigma$ , Lemma 6 then states that for any two linear

transformations  $A$  and  $B$ , the transformed data  $AX$  is to be preferred over the transformed data  $BX$  if

$$(\mu_1 - \mu_0)^\top A^\top (A\Sigma A^\top)^{-1} A(\mu_1 - \mu_0) > (\mu_1 - \mu_0)^\top B^\top (B\Sigma B^\top)^{-1} B(\mu_1 - \mu_0). \quad (33)$$

In particular, using Lemma 6, we obtain the following result showing the dominance of LOL over reduced-rank LDA (or simply `rrLDA` for brevity) when the class conditional distributions are multivariate normal with a common variance.

**Theorem 1.** *Let  $F_0 = N(\mu_0, \Sigma)$  and  $F_1 \sim N(\mu_1, \Sigma)$  be multivariate normal distributions in  $\mathbb{R}^p$ . Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  be the eigenvalues of  $\Sigma$  and  $u_1, u_2, \dots, u_p$  the corresponding eigenvectors. For  $d \leq p$ , let  $U_d = [u_1 \mid u_2 \mid \dots \mid u_d] \in \mathbb{R}^{p \times d}$  be the matrix whose columns are the eigenvectors  $u_1, u_2, \dots, u_d$ . Let  $A = [\delta \mid U_{d-1}]$  and  $B = U_d$  be the LOL and `rrLDA` linear transformations into  $\mathbb{R}^d$ , respectively. Then*

$$\begin{aligned} C(F_0^{(A)}, F_1^{(A)}) - C(F_0^{(B)}, F_1^{(B)}) &= \frac{(\delta^\top (I - U_{d-1} U_{d-1}^\top) \delta)^2}{\delta^\top (\Sigma - \Sigma_{d-1}) \delta} - \delta^\top (\Sigma_d^\dagger - \Sigma_{d-1}^\dagger) \delta \\ &\geq \frac{1}{\lambda_d} \delta^\top (I - U_{d-1} U_{d-1}^\top) \delta - \frac{1}{\lambda_d} \delta^\top (U_d U_d^\top - U_{d-1} U_{d-1}^\top) \delta \geq 0 \end{aligned} \quad (34)$$

and the inequality is strict whenever  $\delta^\top (I - U_d U_d^\top) \delta > 0$ .

*Proof.* We first note that

$$A\Sigma A^\top = [\delta \mid U_{d-1}]^\top \Sigma [\delta \mid U_{d-1}] = \begin{bmatrix} \delta^\top \Sigma \delta & \delta^\top \Sigma U_{d-1} \\ U_{d-1}^\top \Sigma \delta & U_{d-1}^\top \Sigma U_{d-1} \end{bmatrix} = \begin{bmatrix} \delta^\top \Sigma \delta & \delta^\top \Sigma U_{d-1} \\ U_{d-1}^\top \Sigma \delta & \Lambda_{d-1} \end{bmatrix} \quad (35)$$

where  $\Lambda_{d-1} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{d-1})$  is the  $(d-1) \times (d-1)$  diagonal matrix formed by the eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_{d-1}$ . Therefore, letting  $\gamma = \delta^\top \Sigma \delta - \delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta$ , we have

$$\begin{aligned} (A\Sigma A^\top)^{-1} &= \begin{bmatrix} \delta^\top \Sigma \delta & \delta^\top \Sigma U_{d-1} \\ U_{d-1}^\top \Sigma \delta & U_{d-1}^\top \Sigma U_{d-1} \end{bmatrix}^{-1} \\ &= \begin{bmatrix} \gamma^{-1} & -\delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1} \gamma^{-1} \\ -\Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta \gamma^{-1} & (\Lambda_{d-1} - \frac{U_{d-1}^\top \Sigma \delta \delta^\top \Sigma U_{d-1}}{\delta^\top \Sigma \delta})^{-1} \end{bmatrix}. \end{aligned} \quad (36)$$

The Sherman-Morrison-Woodbury formula then implies

$$\begin{aligned} \left( \Lambda_{d-1} - \frac{U_{d-1}^\top \Sigma \delta \delta^\top \Sigma U_{d-1}}{\delta^\top \Sigma \delta} \right)^{-1} &= \Lambda_{d-1}^{-1} + \frac{\Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta \delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1} / (\delta^\top \Sigma \delta)}{1 - \delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta / (\delta^\top \Sigma \delta)} \\ &= \Lambda_{d-1}^{-1} + \frac{\Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta \delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1}}{\delta^\top \Sigma \delta - \delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta} \\ &= \Lambda_{d-1}^{-1} + \gamma^{-1} \Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta \delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1} \end{aligned} \quad (37)$$

We note that  $\Sigma U_{d-1} = U_{d-1} \Lambda_{d-1}$  and  $U_{d-1}^\top \Sigma = \Lambda_{d-1} U_{d-1}^\top$  and hence

$$\begin{aligned} \gamma &= \delta^\top \Sigma \delta - \delta^\top \Sigma U_{d-1} \Lambda_{d-1}^{-1} U_{d-1}^\top \Sigma \delta = \delta^\top \Sigma \delta - \delta^\top U_{d-1} \Lambda_{d-1} \Lambda_{d-1}^{-1} \Lambda_{d-1} U_{d-1}^\top \delta \\ &= \delta^\top \Sigma \delta - \delta^\top U_{d-1} \Lambda_{d-1} U_{d-1}^\top \delta = \delta^\top (\Sigma - \Sigma_{d-1}) \delta \end{aligned} \quad (38)$$

where  $\Sigma_{d-1} = U_{d-1}\Lambda_{d-1}U_{d-1}^\top$  is the best rank  $d - 1$  approximation to  $\Sigma$  with respect to any unitarily invariant norm. In addition,

$$\Lambda_{d-1}^{-1}U_{d-1}^\top\Sigma\delta\delta^\top\Sigma U_{d-1}\Lambda_{d-1}^{-1} = \Lambda_{d-1}^{-1}\Lambda_{d-1}U_{d-1}^\top\delta\delta^\top U_{d-1}\Lambda_{d-1}\Lambda_{d-1}^{-1} = U_{d-1}^\top\delta\delta^\top U_{d-1}.$$

We thus have

$$(A\Sigma A^\top)^{-1} = \begin{bmatrix} \gamma^{-1} & -\delta^\top\Sigma U_{d-1}\Lambda_{d-1}^{-1}\gamma^{-1} \\ -\Lambda_{d-1}^{-1}U_{d-1}^\top\Sigma\delta\gamma^{-1} & (\Lambda_{d-1} - \frac{U_{d-1}^\top\Sigma\delta\delta^\top\Sigma U_{d-1}}{\delta^\top\Sigma\delta})^{-1} \end{bmatrix} \quad (39)$$

$$= \begin{bmatrix} \gamma^{-1} & -\gamma^{-1}\delta^\top U_{d-1} \\ -\gamma^{-1}U_{d-1}^\top\delta & \Lambda_{d-1}^{-1} + \gamma^{-1}U_{d-1}^\top\delta\delta^\top U_{d-1} \end{bmatrix}. \quad (40)$$

Therefore,

$$\begin{aligned} \delta^\top A^\top (A\Sigma A^\top)^{-1} A\delta &= \delta^\top [\delta \mid U_{d-1}] \begin{bmatrix} \gamma^{-1} & -\gamma^{-1}\delta^\top U_{d-1} \\ -\gamma^{-1}U_{d-1}^\top\delta & \Lambda_{d-1}^{-1} + \gamma^{-1}U_{d-1}^\top\delta\delta^\top U_{d-1} \end{bmatrix} [\delta \mid U_{d-1}]^\top \delta \\ &= [\delta^\top \delta \mid \delta^\top U_{d-1}] \begin{bmatrix} \gamma^{-1} & -\gamma^{-1}\delta^\top U_{d-1} \\ -\gamma^{-1}U_{d-1}^\top\delta & \Lambda_{d-1}^{-1} + \gamma^{-1}U_{d-1}^\top\delta\delta^\top U_{d-1} \end{bmatrix} \begin{bmatrix} \delta^\top \delta \\ U_{d-1}^\top \delta \end{bmatrix} \\ &= \gamma^{-1}(\delta^\top \delta)^2 - 2\gamma^{-1}\delta^\top \delta\delta^\top U_{d-1}U_{d-1}^\top\delta + \delta^\top U_{d-1}(\Lambda_{d-1}^{-1} + \gamma^{-1}U_{d-1}^\top\delta\delta^\top U_{d-1})U_{d-1}^\top\delta \\ &= \gamma^{-1}(\delta^\top \delta - \delta^\top U_{d-1}U_{d-1}^\top\delta)^2 + \delta^\top U_{d-1}\Lambda_{d-1}^{-1}U_{d-1}^\top\delta \\ &= \gamma^{-1}(\delta^\top (I - U_{d-1}U_{d-1}^\top)\delta)^2 + \delta^\top \Sigma_{d-1}^\dagger \delta \end{aligned} \quad (41)$$

where  $\Sigma_{d-1}^\dagger$  is the Moore-Penrose pseudo-inverse of  $\Sigma_{d-1}$ . The LDA projection matrix into  $\mathbb{R}^d$  is given by  $B = U_d^\top$  and hence

$$\delta^\top B^\top (B\Sigma B^\top)^{-1} B\delta = \delta^\top U_d\Lambda_d^{-1}U_d^\top\delta = \delta^\top \Sigma_d^\dagger \delta. \quad (42)$$

We thus have

$$\begin{aligned} C(F_0^{(A)}, F_1^{(A)}) - C(F_0^{(B)}, F_1^{(B)}) &= \gamma^{-1}(\delta^\top (I - U_{d-1}U_{d-1}^\top)\delta)^2 - \delta^\top (\Sigma_d^\dagger - \Sigma_{d-1}^\dagger)\delta \\ &= \frac{(\delta^\top (I - U_{d-1}U_{d-1}^\top)\delta)^2}{\delta^\top (\Sigma - \Sigma_{d-1})\delta} - \delta^\top (\Sigma_d^\dagger - \Sigma_{d-1}^\dagger)\delta \\ &\geq \frac{(\delta^\top (I - U_{d-1}U_{d-1}^\top)\delta)^2}{\lambda_d\delta^\top (I - U_{d-1}U_{d-1}^\top)\delta} - \frac{1}{\lambda_d}\delta^\top u_d u_d^\top \delta \\ &= \frac{1}{\lambda_d}\delta^\top (I - U_{d-1}U_{d-1}^\top)\delta - \frac{1}{\lambda_d}\delta^\top (U_d U_d^\top - U_{d-1}U_{d-1}^\top)\delta \geq 0 \end{aligned} \quad (43)$$

where we recall that  $u_d$  is the  $d$ -th column of  $U_d$ . Thus  $C(F_0^{(A)}, F_1^{(A)}) \geq C(F_0^{(B)}, F_1^{(B)})$  always, and the inequality is strict whenever  $\delta^\top (I - U_d U_d^\top)\delta > 0$ .  $\square$

**Remark 1.** *Theorem 1 can be extended to the case wherein the linear transformations are  $A = [\delta \mid U_{d-1}]$  and  $B = U_d$ , respectively, such that  $U_d$  is an arbitrary  $p \times d$  matrix with  $U_d^\top U_d = I$ , and  $U_{d-1}$  is the first  $d - 1$  columns of  $U_d$ . A similar derivation to that in the proof of Theorem 1 then yields*

$$\begin{aligned} C(F_0^{(A)}, F_1^{(A)}) &= \frac{(\delta^\top \Sigma^{-1/2}(I - V_{d-1}V_{d-1}^\top)\Sigma^{1/2}\delta)^2}{\delta^\top \Sigma^{1/2}(I - V_{d-1}V_{d-1}^\top)\Sigma^{1/2}\delta} + \delta^\top \Sigma^{-1/2}V_{d-1}V_{d-1}^\top \Sigma^{-1/2}\delta \\ C(F_0^{(B)}, F_1^{(B)}) &= \delta^\top \Sigma^{-1/2}V_d V_d^\top \Sigma^{-1/2}\delta \end{aligned}$$

where  $V_d V_{d-1}^\top = \Sigma^{1/2} U_d (U_d^\top \Sigma U_d)^{-1} U_d^\top \Sigma^{1/2}$  is the orthogonal projection onto the column space of  $\Sigma^{1/2} U_d$ . Hence  $C(F_0^{(A)}, F_1^{(A)}) > C(F_0^{(B)}, F_1^{(B)})$  if and only if

$$\frac{(\delta^\top \Sigma^{-1/2} (I - V_{d-1} V_{d-1}^\top) \Sigma^{1/2} \delta)^2}{\delta^\top \Sigma^{1/2} (I - V_{d-1} V_{d-1}^\top) \Sigma^{1/2} \delta} > \delta^\top \Sigma^{-1/2} (V_d V_d^\top - V_{d-1} V_{d-1}^\top) \Sigma^{-1/2} \delta. \quad (44)$$

We recover Supp. Eq. 34 by letting  $U_d$  be the matrix whose columns are the eigenvectors corresponding to the  $d$  largest eigenvalue of  $\Sigma$ .

We next present a result relating the Chernoff information for LOL and rRLDA.

**Theorem 2.** Assume the setting of Theorem 1. Let  $C = \tilde{U}_d^\top$  where  $\tilde{U}_d$  is the  $p \times d$  matrix whose columns are the  $d$  largest eigenvectors of the pooled covariance matrix  $\tilde{\Sigma} = \mathbb{E}[(X - \frac{\mu_0 + \mu_1}{2})(X - \frac{\mu_0 + \mu_1}{2})^\top]$ . Then  $C$  is the linear transformation for PCA and

$$\begin{aligned} C(F_0^{(A)}, F_1^{(A)}) - C(F_0^{(C)}, F_1^{(C)}) &= \frac{(\delta^\top (I - U_{d-1} U_{d-1}^\top) \delta)^2}{\delta^\top (\Sigma - \Sigma_{d-1}) \delta} + \delta^\top \Sigma_{d-1}^\dagger \delta - \delta^\top \tilde{\Sigma}_d^\dagger \delta - \frac{(\delta^\top \tilde{\Sigma}_d^\dagger \delta)^2}{4 - \delta^\top \tilde{\Sigma}_d^\dagger \delta} \\ &= \frac{(\delta^\top (I - U_{d-1} U_{d-1}^\top) \delta)^2}{\delta^\top (\Sigma - \Sigma_{d-1}) \delta} + \delta^\top \Sigma_{d-1}^\dagger \delta - \frac{4\delta^\top \tilde{\Sigma}_d^\dagger \delta}{4 - \delta^\top \tilde{\Sigma}_d^\dagger \delta}. \end{aligned} \quad (45)$$

where  $\tilde{\Sigma}_d = \tilde{U}_d \tilde{S}_d \tilde{U}_d^\top$  is the best rank  $d$  approximation to  $\tilde{\Sigma} = \Sigma + \frac{1}{4} \delta \delta^\top$ .

*Proof.* Assume, without loss of generality, that  $\mu_1 = -\mu_0 = \mu$ . We then have

$$\tilde{\Sigma} = \mathbb{E}[X X^\top] = \pi \Sigma + \pi \mu_0 \mu_0^\top + (1 - \pi) \Sigma + (1 - \pi) \mu_1 \mu_1^\top = \Sigma + \mu \mu^\top = \Sigma + \frac{1}{4} \delta \delta^\top. \quad (46)$$

Therefore

$$(C \Sigma C^\top)^{-1} = (\tilde{U}_d^\top \Sigma \tilde{U}_d)^{-1} = (\tilde{U}_d^\top (\tilde{\Sigma} - \frac{1}{4} \delta \delta^\top) \tilde{U}_d)^{-1} = (\tilde{S}_d - \frac{1}{4} \tilde{U}_d^\top \delta \delta^\top \tilde{U}_d)^{-1} \quad (47)$$

$$= \tilde{S}_d^{-1} + \frac{\tilde{S}_d^{-1} \tilde{U}_d^\top \delta \delta^\top \tilde{U}_d \tilde{S}_d^{-1}}{4 - \delta^\top \tilde{U}_d \tilde{S}_d^{-1} \tilde{U}_d^\top \delta} \quad (48)$$

where  $\tilde{S}_d$  is the diagonal matrix containing the  $d$  largest eigenvalues of  $\tilde{\Sigma}$ . Hence

$$\begin{aligned} C(F_0^{(C)}, F_1^{(C)}) &= \delta^\top C^\top (C \Sigma C^\top)^{-1} C \delta = \delta^\top \tilde{U}_d \left( \tilde{S}_d^{-1} + \frac{\tilde{S}_d^{-1} \tilde{U}_d^\top \delta \delta^\top \tilde{U}_d \tilde{S}_d^{-1}}{4 - \delta^\top \tilde{U}_d \tilde{S}_d^{-1} \tilde{U}_d^\top \delta} \right) \tilde{U}_d^\top \delta \\ &= \delta^\top \tilde{U}_d \tilde{S}_d^{-1} \tilde{U}_d^\top \delta + \frac{(\delta^\top \tilde{U}_d \tilde{S}_d^{-1} \tilde{U}_d^\top \delta)^2}{4 - \delta^\top \tilde{U}_d \tilde{S}_d^{-1} \tilde{U}_d^\top \delta} \\ &= \delta^\top \tilde{\Sigma}_d^\dagger \delta + \frac{(\delta^\top \tilde{\Sigma}_d^\dagger \delta)^2}{4 - \delta^\top \tilde{\Sigma}_d^\dagger \delta} = \frac{4\delta^\top \tilde{\Sigma}_d^\dagger \delta}{4 - \delta^\top \tilde{\Sigma}_d^\dagger \delta}. \end{aligned} \quad (49)$$

as desired.  $\square$

**Remark 2.** We recall that the LOL projection  $A = [\delta \mid U_{d-1}]^\top$  yields

$$C(F_0^{(A)}, F_1^{(A)}) = \frac{(\delta^\top (I - U_{d-1} U_{d-1}^\top) \delta)^2}{\delta^\top (\Sigma - \Sigma_{d-1}) \delta} + \delta^\top \Sigma_{d-1}^\dagger \delta. \quad (50)$$

To illustrate the difference between the  $\text{LOL}$  projection and that based on the eigenvectors of the pooled covariance matrix, consider the following simple example. Let  $\Sigma = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$  be a diagonal matrix with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ . Also let  $\delta = (0, 0, \dots, 0, s)$ . Suppose furthermore that  $\lambda_p + s^2/4 < \lambda_d$ . Then we have  $\tilde{\Sigma}_d = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d, 0, 0, \dots, 0)$ . Thus  $\tilde{\Sigma}_d^\dagger = \text{diag}(1/\lambda_1, 1/\lambda_2, \dots, 1/\lambda_d, 0, 0, \dots, 0)$  and  $\delta^\dagger \tilde{\Sigma}_d^\dagger \delta = 0$ . Therefore,  $C(F_0^{(B)}, F_1^{(B)}) = 0$ .

On the other hand, we have

$$C(F_0^{(A)}, F_1^{(A)}) = \frac{(\delta^\top (I - U_{d-1} U_{d-1}^\top) \delta)^2}{\delta^\top (\Sigma - \Sigma_{d-1}) \delta} + \delta^\top \Sigma_{d-1}^\dagger \delta = \frac{s^4}{s^2 \lambda_p} + 0 = s^2 / \lambda_p. \quad (51)$$

A more general form of the previous observation is the following result which shows that  $\text{LOL}$  is preferable over  $\text{PCA}$  when the dimension  $p$  is sufficiently large.

**Proposition 1.** Let  $\Sigma$  be a  $p \times p$  covariance matrix of the form

$$\Sigma = \begin{bmatrix} \Sigma_d & 0 \\ 0 & \Sigma_d^\perp \end{bmatrix} \quad (52)$$

where  $\Sigma_d$  is a  $d \times d$  matrix. Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  be the eigenvalues of  $\Sigma$ , with  $\lambda_1, \lambda_2, \dots, \lambda_d$  being the eigenvalues of  $\Sigma_d$ . Suppose that the entries of  $\delta$  are i.i.d. with the following properties.

1.  $\delta_i \sim Y_i * N(\tau, \sigma^2)$  where  $Y_1, Y_2, \dots, Y_p \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(1 - \epsilon)$ .
2.  $p(1 - \epsilon) \rightarrow \theta$  as  $p \rightarrow \infty$  for some constant  $\theta$ .

Then there exists a constant  $C > 0$  such that if  $\lambda_d - \lambda_{d+1} \geq C\theta\tau^2 \log p$ , then, with probability at least  $\epsilon^d$

$$C(F_0^{(A)}, F_1^{(A)}) > C(F_0^{(B)}, F_1^{(B)}) = 0 \quad (53)$$

*Proof.* The above construction of  $\Sigma$  and  $\delta$  implies, with probability at least  $\epsilon^d$ , that the covariance matrix for  $\tilde{\Sigma}$  is of the form

$$\tilde{\Sigma} = \begin{bmatrix} \Sigma_d & 0 \\ 0 & \Sigma_d^\perp + \frac{1}{4}(\tilde{\delta}\tilde{\delta}^\top) \end{bmatrix} \quad (54)$$

where  $\tilde{\delta} \in \mathbb{R}^{p-d}$  is formed by excluding the first  $d$  elements of  $\delta$ . Now, if  $\lambda_{d+1} + \frac{1}{4}\|\tilde{\delta}\|^2 < \lambda_d$ , then the  $d$  largest eigenvalues of  $\tilde{\Sigma}$  are still  $\lambda_1, \lambda_2, \dots, \lambda_d$ , and thus the eigenvectors corresponding to the  $d$  largest eigenvalues of  $\tilde{\Sigma}$  are the same as those for the  $d$  largest eigenvalues of  $\Sigma$ . That is to say,

$$\lambda_{d+1} + \frac{1}{4}\|\tilde{\delta}\|^2 < \lambda_d \implies \tilde{\Sigma}_d^\dagger = \Sigma_d^\dagger \implies \delta^\top \tilde{\Sigma}_d^\dagger \delta = 0 \implies C(F_0^{(B)}, F_1^{(B)}) = 0. \quad (55)$$

We now compute the probability that  $\lambda_{d+1} + \frac{1}{4}\|\tilde{\delta}\|^2 < \lambda_d$ . Suppose for now that  $\epsilon > 0$  is fixed and does not vary with  $p$ . We then have

$$\frac{\sum_{i=d+1}^p \delta_i^2 - (p-d)(1-\epsilon)\tau^2}{\sqrt{(p-d)(2(1-\epsilon)(2\tau^2\sigma^2 + \sigma^4) + \epsilon(1-\epsilon)(\tau^4 + 2\tau^2\sigma^2 + \sigma^4))}} \xrightarrow{d} N(0, 1). \quad (56)$$

Thus, as  $p \rightarrow \infty$ , the probability that  $\lambda_{d+1} + \frac{1}{4}\|\tilde{\delta}\|^2 < \lambda_d$  converges to that of

$$\Phi\left(\frac{4(\lambda_d - \lambda_{d+1}) - (p-d)(1-\epsilon)\tau^2}{\sqrt{(p-d)(2(1-\epsilon)(2\tau^2\sigma^2 + \sigma^4) + \epsilon(1-\epsilon)(\tau^4 + 2\tau^2\sigma^2 + \sigma^4))}}\right). \quad (57)$$

This probability can be made arbitrarily close to 1 provided that  $\lambda_d - \lambda_{d+1} \geq Cp(1 - \epsilon)\tau^2$  for all sufficiently large  $p$  and for some constant  $C > 1/4$ . Since the probability that  $\delta_1 = \delta_2 = \dots = \delta_d$  is at least  $\epsilon^d$ , we thus conclude that for sufficiently large  $p$ , with probability at least  $\epsilon^d$ ,

$$C(F_0^{(B)}, F_1^{(B)}) = 0 < C(F_0^{(A)}, F_1^{(A)}). \quad (58)$$

In the case where  $\epsilon = \epsilon(p) \rightarrow 1$  as  $p \rightarrow \infty$  such that  $p(1 - \epsilon) \rightarrow \theta$  for some constant  $\theta$ , then the probability that  $\lambda_{d+1} + \frac{1}{4}\|\delta\|^2 < \lambda_d$  converges to the probability that

$$\frac{1}{4} \sum_{i=1}^K \sigma^2 \chi_1^2(\tau) \geq \lambda_d - \lambda_{d+1} \quad (59)$$

where  $K$  is Poisson distributed with mean  $\theta$  and  $\chi_i^2(\tau)$  is the non-central chi-square distribution with one degree of freedom and non-centrality parameter  $\tau$ . Thus if  $\lambda_d - \lambda_{d+1} \geq C\theta\tau^2 \log p$  for sufficiently large  $p$  and for some constant  $C$ , then this probability can also be made arbitrarily close to 1.  $\square$

**Remark 3.** *The previous comparisons are done for the case of  $C = 2$  classes. Extending these comparisons to the case of  $C > 2$  classes is, however, non-trivial. More precisely, suppose we have  $Y \in \{1, 2, \dots, C\}$  and that, conditional on  $Y = c$ ,  $X \sim \mathcal{N}(\mu_c, \Sigma)$  is multivariate normal with mean  $\mu_c$  and common covariance matrix  $\Sigma$ . Then, given  $X = x$ , the Bayes optimal classifier for  $Y$  is still*

$$g_{\text{LDA}}(x) = \operatorname{argmin}_{y \in \{1, 2, \dots, C\}} \left[ \frac{1}{2} (x - \mu_y)^\top \Sigma^{-1} (x - \mu_y) - \log \pi_y \right] \quad (60)$$

$$= \operatorname{argmin}_{y \in \{1, 2, \dots, C\}} \left[ -x^\top \Sigma^{-1} \mu_y + \frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y - \log \pi_y \right] \quad (61)$$

Taking  $\frac{1}{2} \mu_y^\top \Sigma^{-1} \mu_y - \log \pi_y$  as either a given constant or as an intercept term to be learned or estimated, the reduced-rank LDA for  $C > 2$  classes still corresponds to looking at the top  $d$  eigenvectors of  $\Sigma$ . That is to say, we transform the predictor variables via  $x \mapsto U_d x$  followed by performing LDA on the transformed data. Similarly, the PCA transformation corresponds to using the top  $d$  eigenvectors of the pooled covariance matrix  $\tilde{\Sigma} = \mathbb{E}[(X - \sum_c \pi_c \mu_c)(X - \sum_c \pi_c \mu_c)^\top]$  followed by performing LDA. Suppose we now compare LOL, rrLDA, and PCA in this multi-class setting. Let  $A: X \mapsto AX$  be a linear transformation. Then by Supp. Eq. (30) and Supp. Eq. (32), the Chernoff information for the transformed data in this multi-class setting is

$$\min_{c \neq c'} \frac{1}{8} (\mu_c - \mu_{c'})^\top A^\top (A \Sigma A^\top)^{-1} A (\mu_c - \mu_{c'}). \quad (62)$$

We now see that, in the case of rrLDA and PCA the linear transformation  $A$  depends only on the covariance matrix  $\Sigma$  and  $\tilde{\Sigma}$ , respectively. That is to say, the linear transformation  $A$  does not depend on the choice of  $c$  and  $c'$ . In contrast, currently for LOL the linear transformation  $A$  depends on both  $\Sigma$  as well as  $\mu_c - \mu_{c'}$ . In other words, there is no single choice for  $A$  but rather that  $A$  changes as  $c, c'$  changes. Direct comparison, in the multi-classes setting, between LOL and either of rrLDA or PCA is thus an open problem that we leave for future work. Finally we note that if we allow the linear transformation for LOL to vary with the classes  $c$  and  $c'$ , i.e., taking a one-vs-one approach to multi-classes classification, then the results presented in this paper are valid for all pairs  $c, c'$ .



### 4.3 Finite Sample Performance

We now consider the finite sample performance of LOL and PCA-based classifiers in the high-dimensional setting with small or moderate sample sizes, e.g., when  $p$  is comparable to  $n$  or when  $n \ll p$ . Once again we assume that  $X|Y = i \sim \mathcal{N}(\mu_i, \Sigma)$  for  $i = 0, 1$ . Furthermore, we also assume that  $\Sigma$  belongs to the class  $\Theta(p, r, k, \tau, \lambda)$  as defined below.

**Definition** Let  $\lambda > 0$ ,  $\tau \geq 1$  and  $k \leq p$  be given. Denote by  $\Theta(p, r, k, \tau, \lambda, \sigma^2)$  the collection of matrices  $\Sigma$  such that

$$\Sigma = V\Lambda V^\top + \sigma^2 I \quad (63)$$

where  $V$  is a  $p \times r$  matrix with orthonormal columns and  $\Lambda$  is a  $r \times r$  diagonal matrix whose diagonal entries  $\lambda_1, \lambda_2, \dots, \lambda_r$  satisfy  $\lambda \geq \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r \geq \lambda/\tau$ . In addition, assume also that  $|\text{supp}(V)| \leq k$  where  $\text{supp}(V)$  denote the non-zero rows of  $V$ , i.e.,  $\text{supp}(V)$  is the subset of  $\{1, 2, \dots, p\}$  such that  $V_j \neq 0$  if and only if  $j \in \text{supp}(V)$ .

We note that in general  $r \leq k \ll p$  and  $\lambda/\tau \gg \sigma^2$ . We then have the following result.

**Theorem 3** ([4]). *Suppose there exist constants  $M_0$  and  $M_1$  such that  $M_1 \log p \geq \log n \geq M_0 \log \lambda$ . Then there exists a constant  $c_0 = c_0(M_0, M_1)$  depending on  $M_0$  and  $M_1$  such that for all  $n$  and  $p$  for which*

$$\frac{\tau k}{n} \log \frac{ep}{k} \leq c_0, \quad (64)$$

there exists an estimate  $\hat{V}$  of  $V$  such that

$$\sup_{\Sigma \in \Theta(p, r, k, \tau, \lambda, \sigma^2)} \mathbb{E} \|\hat{V}\hat{V}^\top - VV^\top\|^2 \leq \frac{Ck(\sigma\lambda + \sigma^2)}{n\lambda^2} \log \frac{ep}{k} \quad (65)$$

where  $C$  is a universal constant not depending on  $p, r, k, \tau, \lambda$  and  $\sigma^2$ .

Theorem 3 then implies the following result for comparing the Chernoff information of the sample version of LOL against that for PCA.

**Corollary 4.** *Let  $\Sigma \in \Theta(p, r, k, \tau, \lambda)$  as defined above. Suppose that  $C(F_0^{(A)}, F_1^{(A)}) > C(F_0^{(B)}, F_1^{(B)})$  where  $A$  and  $B$  denote the LOL and PCA projection matrices based on the eigenvectors of  $\Sigma$  associated with the  $d \leq r$  largest eigenvalues, i.e.,  $A = [\delta | V_{1:d-1}]$  and  $B = V_{1:d}$ . Then there exists constants  $M$  and  $c$  such that if  $\log n \geq M \log \lambda$  and  $\frac{\tau k}{n} \log \frac{ep}{k} \leq c$ , then there exists an estimate  $\hat{V}$  of  $V$  such that, with  $\hat{A} = [\hat{\delta} | \hat{V}_{1:d-1}]$  and  $\hat{B} = [\hat{V}_{1:d}]$ , we have*

$$\mathbb{E}[C(F_0^{(\hat{A})}, F_1^{(\hat{A})})] > \mathbb{E}[C(F_0^{(\hat{B})}, F_1^{(\hat{B})})] \quad (66)$$

The above corollary states that for  $\Sigma \in \Theta(p, r, k, \tau, \lambda)$ , then provided that the Chernoff information of the population version of LOL is larger than the Chernoff information of the population version of PCA, we can choose  $n$  sufficiently large (as compared to  $\lambda$  and  $\tau$  and  $k$ ) such that the expected Chernoff information for the sample version of LOL is also larger than the expected Chernoff information of the sample version of PCA. We emphasize that it is necessary that the LOL and the PCA version are both

projected into the top  $d \leq r$  dimension of the sample covariance matrices. The constants  $M$  and  $c$  in the statement of the above corollary are chosen so that  $M$  (which depends on  $M_0$  and  $M_1$  in the statement of Theorem 3) is sufficiently large and  $c$  (which depends on  $c_0$ ) is sufficiently small to ensure that the bound in Supp. Eq. (65) is sufficiently small. If  $C(F_0^{(A)}, F_1^{(A)}) > C(F_0^{(B)}, F_1^{(B)})$  and  $\|\hat{V}\hat{V}^\top - VV^\top\|$  is sufficiently small, then  $\mathbb{E}[C(F_0^{(\hat{A})}, F_1^{(\hat{A})})] > \mathbb{E}[C(F_0^{(\hat{B})}, F_1^{(\hat{B})})]$  as desired.

## Supplementary Note 5 Real-Data Performance Analysis

PCA, the industry-standard dimensionality reduction technique for high-dimensional problems, is compared to LOL, RP, and rRLDA in terms of cross-validated classification error.

In the experiments, we used  $k$  fold cross-validation. Testing sets were rotated across all folds, with the training sets comprising the remaining  $k - 1$  folds. A low-dimensional projection matrix  $\mathbf{A}$  is first learned through the training set, and the low-dimensional training points are then used to train an LDA classifier  $C$ . The testing points are embedded via  $\mathbf{A}$ , and classification error is determined using the trained classifier  $C$ .

Performance is assessed using Cohen’s kappa [7], which normalizes the classification error typically between 0 (the classifier performs no better than the classifier which guesses the most-likely class from the training set, the `chance` classifier) and 1 (the trained classifier performs perfectly). Negative scores can be achieved if the trained classifier performs worse than the `chance` classifier. The effect size is measured as the difference between Cohen’s kappa for the trained classifier after embedding with  $\text{PCA}\kappa(\text{PCA})$  and the trained classifier after embedding with technique  $\varepsilon$ ,  $\kappa(\varepsilon)$ . Table 1 provides details about each neuroimaging dataset.

Problem	Sample Size ( $n$ )	Training Size*	# Features ( $p$ )	Classes ( $K$ )	Source
Templeton114	111	100	$> 1.5 \times 10^8$	2	MRN
BNU1	110	100	$> 1.5 \times 10^8$	2	CoRR [22]
BNU3	47	43	$> 1.5 \times 10^8$	2	CoRR [22]
SWU4	453	407	$> 1.5 \times 10^8$	2	CoRR [22]
KKI2009	42	38	$> 1.5 \times 10^8$	2	KKI [12]
Genomics	340	306	745,184	2	Douville et al. [9]

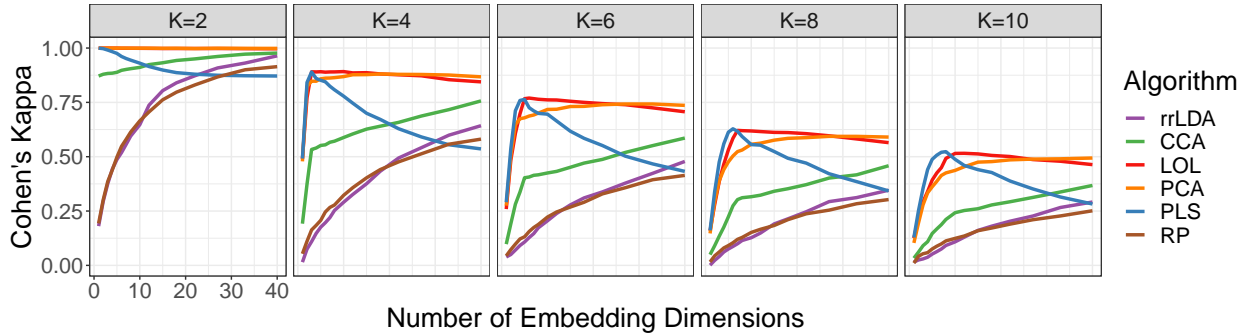
**Supplementary Table 1.** Table of datasets used in this study. The top 5 datasets (neuroimaging) are pre-processing by only registering the brains to the MNI152 template[11, 17, 19, 20]. The neuroimaging dataset comprises a total of 5 classification problems (5 datasets across a single sex classification task). The bottom dataset (genomics) is pre-processed by aligning sequencing data to 745,184 amplicons on the human genome. The genomics dataset comprises two benchmark classification problems (sex or age).

## Supplementary Note 6 Extensions to Other Supervised Learning Problems

### 6.1 Large numbers of classes

Here, we explore an experiment in which the number of classes increases for a given simulation. We look at the multiclass hump- $K$  problem, described in Section Supplementary Note 3. In this simulation, while the space spanned by the differences of means conveys more information than the directions of

## Performance of Embedding Strategies on Multiclass Problems

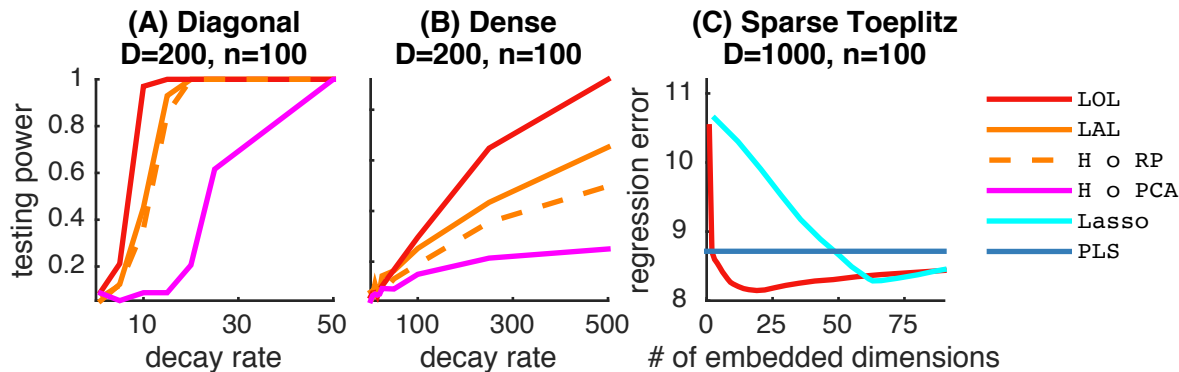


**Supplementary Figure 1. The Multiclass Hump Simulation.** We show the results of the multiclass trunk problem, as the number of classes increases from 2 to 10, with the number of dimensions and the number of samples fixed. Effect size is measured with Cohen's Kappa. `LOL` and `PLS` provide better performance over competing techniques including `PCA`, and this gap widens as the number of classes increases.

maximal variance, we expect that the shift in the means for a given class at a given dimension should also increase the variance fractionally in that direction as well. Supplementary Figure 1A shows the simulation setup, for  $K = 10$ . Supplementary Figure 1B indicates the misclassification rate as a function of Cohen's Kappa. We use Cohen's Kappa instead of the misclassification rate for direct evaluation since  $K$  varies widely across these simulations at a fixed number of total samples  $n = 128$ , making the difficulty of the problem as  $K$  increases two-fold: not only are there more classes, but there are also fewer examples of each class per simulation setting. In all cases, the best random classifier would be the classifier that continually guesses a single class continuously, which has expected accuracy of  $\frac{1}{K}$ . On all simulations, we see that both `PLS` and `LOL` rapidly approach a higher Kappa statistic (better performance relative the random classifier) as they learn the space spanned by the differences of means. `PLS` rapidly declines in performance as successive dimensions are added, and `LOL` sees a small performance decline, as successive dimensions should convey no information regarding the class. `PCA` is able to ultimately identify the space spanned by the differences of the means, but takes far more embedding dimensions to do so, and yields a lower Kappa statistic than either of the other two strategies.

## 6.2 Hypothesis Testing

The utility of incorporating the mean difference vector into supervised machine learning extends beyond classification. In particular, hypothesis testing can be considered as a special case of classification, with a particular loss function. We therefore apply the same idea to a hypothesis testing scenario. The multivariate generalization of the t-test, called Hotelling's Test, suffers from the same problem as does the classification problem; namely, it requires inverting an estimate of the covariance matrix, which would result in a matrix that is low-rank and therefore singular in the high-dimensional setting. To mitigate this issue in the hypothesis testing scenario, prior work applied similar tricks as they have done in the classification setting. One particularly nice and related example is that of Lopes et al. [14], who addresses this dilemma by using random projections to obtain a low-dimensional representation, following by applying Hotelling's Test in the lower-dimensional subspace. Supplementary Figure 2A and B show the power of their test (labeled `RP`) alongside the power of `PCA`, `LOL`, and `LFL` for two different conditions. In each case we use the different approaches to project to low dimensions, followed by using Hotelling's test on the projected data. In the first example the true covariance matrix is diagonal, and in



**Supplementary Figure 2.** The intuition of including the mean difference vector is equally useful for other supervised manifold learning problems, including testing and regression. **(A)** and **(B)** show two different high-dimensional testing settings, as described in Methods. Power is plotted against the decay rate of the spectrum, which approximates the effective number of dimensions. LOL composed with Hotelling’s test outperforms the random projections variants described in [14], as well as several other variants. **(C)** A sparse high-dimensional regression setting, as described in Methods, designed for sparse methods to perform well. Log<sub>10</sub> mean squared error is plotted against the number of projected dimensions. LOL composed with linear regression outperforms Lasso (cyan), the classic sparse regression method, as well as partial least squares (PLS; black). These three simulation settings therefore demonstrate the generality of this technique.

the second, the true covariance matrix is dense. The horizontal axis on both panels characterizes the decay rate of the eigenvalues, so larger numbers imply the data is closer to low-rank (see Methods for details). The results indicate that the LOL test has higher power for essentially all scenarios. Moreover, it is not merely replacing random projections with PCA (solid magenta line), nor simply incorporating the mean difference vector (dashed green line), but rather, it appears that LOL for testing uses both modifications to improve performance.

### 6.3 Regression

High-dimensional regression is another supervised learning method that can use the LOL idea. Linear regression, like classification and Hotelling’s Test, requires inverting a matrix as well. By projecting the data onto a lower-dimensional subspace first, followed by linear regression on the low-dimensional data, we can mitigate the curse of high-dimensions. To choose the projection matrix, we partition the data into  $K$  partitions (we select  $K = 10$  arbitrarily), based on the percentile of the target variable, we obtain a  $K$ -class classification problem. Then, we can apply LOL to learn the projection. Supplementary Figure 2C shows an example of this approach, contrasted with Lasso and partial least squares, in a sparse simulation setting (see Methods for details). LOL is able to find a better low-dimensional projection than Lasso, and performs significantly better than partial least squares, for essentially all choices of number of dimensions to project into.

### Supplementary Note 7 The R implementation of LOL

Supplementary Figure 3 shows the R implementation of LOL for binary classification using FlashMatrix [21]. The implementation takes a  $D \times I$  matrix, where each column is a training instance and each instance has  $D$  features, and outputs a  $D \times k$  projection matrix.

```

LOL <- function(m, labels, k) {
  counts <- fm.table(labels)
  num.labels <- length(counts$val)
  num.features <- dim(m)[1]
  nv <- k - (num.labels - 1)
  gr.sum <- fm.groupby(m, 1, fm.as.factor(labels, 2), fm.bo.add)
  gr.mean <- fm.mapply.row(gr.sum, counts$Freq, fm.bo.div, FALSE)
  diff <- fm.get.cols(gr.mean, 1) - fm.get.cols(gr.mean, 2)
  svd <- fm.svd(m, nv=0, nu=nv)
  fm.cbind(diff, svd$u)
}

```

**Supplementary Figure 3.** The R implementation of LOL.

---

**Pseudocode 1** Simple pseudocode for two class LOL on sample data.

---

**Input:**  $X$  a  $p \times n$  matrix ( $n \ll p$ ), where columns are observations; rows are features. An  $n$  length vector of observation labels,  $y$ . An integer  $k$  to specify desired output dimension.

**Output:**  $A \in \mathbb{R}^{p \times k}$

```

1: function LOL.TRAIN( $X, Y, k$ )
2:   for all  $j \in J$  do
3:      $n_j = \sum_{i=1}^n \mathbf{I}(y_i = j)$  ▷ sample size per class
4:      $\hat{\mu}_j = \frac{1}{n_j} \sum_{i=1}^n \mathbf{x}_i \mathbf{I}(y_i = j)$  ▷ class means
5:   end for
6:    $\hat{\delta} = \hat{\mu}_1 - \hat{\mu}_2$  ▷ difference of means
7:    $\hat{\delta} = \hat{\delta} / \|\hat{\delta}\|$  ▷ unit normalize difference of means
8:   for all  $i \in [n]$  do
9:      $\tilde{x}_i = \mathbf{x}_i - \hat{\mu}_{y_i}$  ▷ class centered data
10:  end for
11:   $[\hat{u}, \hat{d}, \hat{v}] = \text{svds}(\tilde{x}, k - 1)$  ▷ compute top  $k$  singular vectors
12:   $A = [\hat{\delta}, \hat{u}]$  ▷ concatenate difference of the means and the top  $k$  right singular vectors
13: end function

```

---

## Supplementary References

- [1] J. Ahn and J. S. Marron. The maximum data piling direction for discrimination. *Biometrika*, 97: 254–259, 2010. **7**
- [2] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammerling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide: Third Edition*. SIAM, 1999. ISBN 0898714478. URL <https://books.google.com/books?hl=en&lr=&id=AZlvEnr9gCgC&pgis=1>. **4**
- [3] L. Breiman. Statistical modeling: The two cultures. *Statistical Science*, 16(3):199–231, 2001. **5**
- [4] Tony Cai, Zongming Ma, and Yihong Wu. Optimal estimation and rank detection for sparse spiked covariance matrices. *Probab. Theory Related Fields*, 161(3-4):781–815, apr 2015. **17**
- [5] William R. Carson, Minhua Chen, Miguel R. D. Rodrigues, Robert Calderbank, and Lawrence Carin. Communications-Inspired Projection Design with Application to Compressive Sensing. *arXiv*, Jun 2012. doi: 10.1137/120878380. **2**
- [6] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Annals of Mathematical Statistics*, 23:493–507, 1952. **10**
- [7] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46, 1960. doi: 10.1177/001316446002000104. URL <https://doi.org/10.1177/001316446002000104>. **18**
- [8] I. Csizár. Information-type measures of difference of probability distributions and indirect observations. *Studia Scientiarum Mathematicarum Hungarica*, 2:229–318, 1967. **11**
- [9] Christopher Douville, Joshua D. Cohen, Janine Ptak, Maria Popoli, Joy Schaefer, Natalie Silliman, Lisa Dobbyn, Robert E. Schoen, Jeanne Tie, Peter Gibbs, Michael Goggins, Christopher L. Wolfgang, Tian-Li Wang, Ie-Ming Shih, Rachel Karchin, Anne Marie Lennon, Ralph H. Hruban, Cristian Tomasetti, Chetan Bettgowda, Kenneth W. Kinzler, Nickolas Papadopoulos, and Bert Vogelstein. Assessing aneuploidy with repetitive element sequencing. *Proc. Natl. Acad. Sci. U.S.A.*, 117(9): 4858–4863, Mar 2020. ISSN 0027-8424. doi: 10.1073/pnas.1910041117. **18**
- [10] Jianqing Fan, Yang Feng, and Xin Tong. A road to classification in high dimensional space: the regularized optimal affine discriminant. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 74(4):745–771, September 2012. ISSN 13697412. doi: 10.1111/j.1467-9868.2012.01029.x. URL <http://doi.wiley.com/10.1111/j.1467-9868.2012.01029.x>. **2**
- [11] Mark Jenkinson et al. FSL. *NeuroImage*, 62(2):782–90, aug 2012. ISSN 1095-9572. URL <http://www.ncbi.nlm.nih.gov/pubmed/21979382>. **18**
- [12] Bennett A Landman, Alan J Huang, Aliya Gifford, Deepti S Vikram, Issel Anne L Lim, Jonathan A D Farrell, John A Bogovic, Jun Hua, Min Chen, Samson Jarso, Seth A Smith, Suresh Joel, Susumu Mori, James J Pekar, Peter B Barker, Jerry L Prince, and Peter C M van Zijl. Multi-parametric neuroimaging reproducibility: a 3-T resource study. *Neuroimage*, 54(4):2854–2866, February 2011. **18**
- [13] C. C. Leang and D. H. Johnson. On the asymptotics of M-hypothesis bayesian detection. *IEEE Transactions on Information Theory*, 43:280–282, 1997. **11**

- [14] Miles Lopes, Laurent Jacob, and Martin J. Wainwright. A More Powerful Two-Sample Test in High Dimensions using Random Projection. In *Neural Information Processing Systems*, pages 1206–1214, 2011. URL <http://papers.nips.cc/paper/4260-a-more-powerful-two-sample-test-in-high-dimensions-using-random-projection>. 19, 20
- [15] Qing Mai, Hui Zou, and Ming Yuan. A direct approach to sparse discriminant analysis in ultra-high dimensions. *Biometrika*, 99(1), Feb 2012. doi: 10.2307/41720670. 2
- [16] K. V. Mardia, J. T. Kent, and J. M. Bibby. *Multivariate Analysis*. Academic Press, 1979. 7
- [17] John Mazziotta et al. A four-dimensional probabilistic atlas of the human brain. *Journal of the American Medical Informatics Association*, 8(5):401–430, 2001. 18
- [18] H. Shin and R. L. Eubank. Unit canonical correlations and high-dimensional discriminant analysis. *Journal of Statistical Computation and Simulation*, 81:167–178, 2011. 7
- [19] Stephen M Smith et al. Advances in functional and structural MR image analysis and implementation as FSL. *NeuroImage*, 23 Suppl 1:S208–19, jan 2004. ISSN 1053-8119. URL <http://www.ncbi.nlm.nih.gov/pubmed/15501092>. 18
- [20] Mark W Woolrich et al. Bayesian analysis of neuroimaging data in FSL. *NeuroImage*, 45(1 Suppl): S173–86, mar 2009. ISSN 1095-9572. URL <http://www.sciencedirect.com/science/article/pii/S1053811908012044>. 18
- [21] Da Zheng, Disa Mhembere, Joshua T Vogelstein, Carey E Priebe, and Randal Burns. Flashmatrix: Parallel, scalable data analysis with generalized matrix operations using commodity ssds. *arXiv preprint arXiv:1604.06414*, 2016. 20
- [22] Xi-Nian Zuo, Jeffrey S Anderson, Pierre Bellec, Rasmus M Birn, Bharat B Biswal, Janusch Blautzik, John CS Breitner, Randy L Buckner, Vince D Calhoun, F Xavier Castellanos, et al. An open science resource for establishing reliability and reproducibility in functional connectomics. *Scientific data*, 1:140049, 2014. 18