# DPCA: DIMENSIONALITY REDUCTION FOR DISCRIMINATIVE ANALYTICS OF MULTIPLE LARGE-SCALE DATASETS

Gang Wang, Jia Chen, and Georgios B. Giannakis

ECE Dept. and Digital Tech. Center, Univ. of Minnesota, Mpls., MN 55455, USA Emails: {gangwang, chen5625, georgios}@umn.edu

# ABSTRACT

Principal component analysis (PCA) has well-documented merits for data extraction and dimensionality reduction. PCA deals with a single dataset at a time, and it is challenged when it comes to analyzing multiple datasets. Yet in certain setups, one wishes to extract the most significant information of one dataset relative to other datasets. Specifically, the interest may be on identifying or extracting features that are specific to a single target dataset but not the others. This paper presents a novel approach for such so-termed discriminative data analysis, and establishes its optimality in the least-squares sense under suitable assumptions. The criterion reveals linear combinations of variables by maximizing the ratio of the variance of the target data to that of the remainders. The novel approach solves a generalized eigenvalue problem by performing SVD just once. Numerical tests using synthetic and real datasets showcase the merits of the proposed approach relative to its competing alternatives.

*Index Terms*— Dimensionality reduction, robust principal component analysis, discriminative analytics

# 1. INTRODUCTION

Principal component analysis (PCA) is arguably the most widely used method for data visualization and dimensionality reduction [1]. PCA originated in statistics [2], but its modern instantiation as well as the term *principal component* (PC) vector was formalized in [3]. The goal of PCA is to extract the most important information from a data table representing observations, and depict it as a few PCs. PCs are uncorrelated linear transformations of the original set of variables, along which the maximum variation in the data is captured [1].

Yet, several application domains involve *multiple* datasets, and the task is to extract trends or features depicted by *component vectors* that are present in one dataset *but not* the other(s) [4]. For example, consider two gene expression datasets of individuals from across different countries and genders: the first includes gene expression levels of cancer patients, which constitutes the *target data* that we want to analyze, while the

second is formed by healthy volunteers, and is called *control* or *background data*. Applying PCA on either the target data or the target augmented with the control data is likely to obtain PCs that represent the background information common to both datasets (e.g., the demographic patterns, genders) [5], instead of the component vectors depicting the subtypes of cancer within patients. Despite its practical relevance, such discriminative data analysis has not been widely studied.

Generalizations to PCA include multi-dimensional scaling [6], local linear embedding [7], robust or kernel PCA [8, 9, 10], and canonical component analysis [11, 12]. Given multiple datasets, analysts have to perform these procedures on each individual dataset, and subsequently evaluate manually the obtained projections to identify whether significant patterns representing similarities or differences across datasets are present. A recent approach pursued what is termed contrastive (c) PCA for extracting the most distinct features of one dataset relative to the other [4]. cPCA is able to reveal the dataset-specific information often missed by PCA [4]. This becomes possible through a hyper-parameter that takes into account both target and control data, and critically affects the performance of cPCA. The cPCA solution is often found with SVD. Although possible to automatically select the best from a list of given values, computing SVD multiple times can be computationally cumbersome or even prohibitive in large-scale data extraction settings. Another method related to PCA is linear discriminant analysis (LDA) [13], that is "supervised," and seeks linear combinations of variables to maximize the separation between classes. This is achieved by maximizing the ratio of the variance between classes to the variance within the classes.

Inspired by LDA and cPCA, this paper puts forth a new method for discriminative analytics, which is shown to be optimal in the least-squares (LS) sense provided that the background component vectors are also present in the target data. Our method seeks linear combinations of variables by maximizing the ratio of the variance of target data to the variance of control data, which justifies our chosen description as *discriminative (d) PCA*. dPCA solves a generalized eigenvalue problem. Relative to cPCA, dPCA is parameter-free, and requires only one SVD. As such, dPCA is well-suited for largescale either discriminative or contrasting data exploration.

Work in this paper was supported in part by NIH 1R01GM104975-01 and NSF 1500713.

# 2. PRELIMINARIES AND PRIOR ART

Consider two datasets, the target data  $\{x_i \in \mathbb{R}^D\}_{1 \le i \le m}$  that we are interested in analyzing, and data  $\{y_j \in \mathbb{R}^D\}_{1 \le j \le n}$ containing latent background component vectors in the target data. Assume without loss of generality (wlog) that the sample mean of each dataset has been removed, and let  $C_{xx} :=$  $(1/m) \sum_{i=1}^m x_i x_i^\top$  and  $C_{yy} := (1/n) \sum_{i=1}^n y_i y_i^\top$  denote the corresponding sample covariance matrices. To motivate the novel approach in Sec. 3, the basics of PCA and cPCA are briefly reviewed in this section.

One formulation of PCA seeks vectors  $\{\chi_i \in \mathbb{R}^d\}_{1 \le i \le m}$ as linear combinations of  $\{x_i \in \mathbb{R}^D\}_{1 \le i \le m}$  with d < D via maximizing their variances in the low-dimensional subspaces [1]. Specifically for d = 1, (linear) PCA obtains  $\chi_i := u^{\top} x_i$ , where the direction  $u \in \mathbb{R}^D$  is found by

$$\max_{\boldsymbol{u}\in\mathbb{R}^D} \quad \boldsymbol{u}^\top \boldsymbol{C}_{xx}\boldsymbol{u} \tag{1a}$$

s. to 
$$\boldsymbol{u}^{\top}\boldsymbol{u} = 1.$$
 (1b)

Solving (1) yields u as the principal eigenvector of matrix  $C_{xx}$ , also known as the first PC. Instead of having constraint (1b) explicitly, we assume wlog that the solution u will always be normalized to have unity norm. For d > 1, PCA amounts to computing the first d principal eigenvectors of  $C_{xx}$  instead. As alluded to in Sec. 1, when two datasets  $\{x_i\}$  and  $\{y_j\}$  are available, PCA performed either on  $\{x_i\}$ , or on  $\{\{x_i\}, \{y_j\}\}$ , can generally not unveil the patterns or trends that are specific to the target relative to the control data.

Contrastive (c) PCA [4], on the other hand, aims to identify directions u along which the target data possesses large variations while the control data has small variations. Concretely, cPCA pursues problem [4]

$$\max_{\|\boldsymbol{u}\|_{2}=1} \boldsymbol{u}^{\top} \boldsymbol{C}_{xx} \boldsymbol{u} - \alpha \boldsymbol{u}^{\top} \boldsymbol{C}_{yy} \boldsymbol{u}.$$
(2)

whose solution is given by the eigenvector of  $C_{xx} - \alpha C_{yy}$  associated with the largest eigenvalue, and constitutes the first contrastive (c) PC. Here,  $\alpha > 0$  is a hyper-parameter that trades off maximizing the target data variance (the first term in (2)) for minimizing the control data variance (second term). However, there is no rule of thumb for choosing  $\alpha$ . Although a spectral clustering based algorithm has been developed to automatically select the value of  $\alpha$ , its brute-force search discourages its use in large-scale datasets.

#### **3. THE NOVEL APPROACH**

Unlike PCA, LDA is "supervised," and seeks directions that yield the largest possible separation between classes via maximizing the ratio of the variance across classes to the variance within classes. In the same vein, when both the target and the background data are available, and one is interested in extracting features, namely component vectors that are *only* present in the target data but *not* in the background data, a meaningful approach would be to maximize the ratio of the variance of the target data over that of the background data

$$\max_{\boldsymbol{u}\|_{2}=1} \quad \frac{\boldsymbol{u}^{\top} \boldsymbol{C}_{xx} \boldsymbol{u}}{\boldsymbol{u}^{\top} \boldsymbol{C}_{yy} \boldsymbol{u}}$$
(3)

which, with slight abuse of the term "discriminant," we call *discriminative (d) PCA*. Likewise, the solution of (3) will be termed first discriminative PC, or dPC for short.

Ш

## 3.1. dPCA Algorithm

Suppose that  $C_{yy}$  is full rank with eigen-decomposition  $C_{yy} := U_y^{\top} \Sigma_y U_y$ . Upon defining  $C_{yy}^{1/2} := \Sigma_y^{1/2} U_y$ , and changing variables  $v := C_{yy}^{1/2} u$ , (3) admits the same solution as

$$\boldsymbol{v}^* := \arg \max_{\|\boldsymbol{v}\|_2 = 1} \quad \boldsymbol{v}^\top \boldsymbol{C}_{yy}^{-\top/2} \boldsymbol{C}_{xx} \boldsymbol{C}_{yy}^{-1/2} \boldsymbol{v}$$
(4)

which is the principal eigenvector of  $C_{yy}^{-\top/2}C_{xx}C_{yy}^{-1/2}$ . Finally, the solution to (3) is recovered as  $u^* := C_{yy}^{-1/2}v^*$ , followed by normalization to obtain a unit norm.

On the other hand, leveraging Lagrangian duality, the solution of (3) can also be obtained as the right eigenvector of  $C_{uy}^{-1}C_{xx}$ . To see this, note that (3) can be rewritten as

$$\max_{\boldsymbol{u}\in\mathbb{R}^D} \quad \boldsymbol{u}^\top \boldsymbol{C}_{xx} \boldsymbol{u} \tag{5a}$$

s. to 
$$\boldsymbol{u}^{\top} \boldsymbol{C}_{yy} \boldsymbol{u} = b$$
 (5b)

for some constant b > 0 such that the solution  $||u^*||_2 = 1$ . One can set b = 1 and subsequently normalize the solution of (5). Letting  $\lambda \in \mathbb{R}$  be the dual variable corresponding to constraint (5b), the Lagrangian of (5) is

$$\mathcal{L}(\boldsymbol{u};\boldsymbol{\lambda}) = \boldsymbol{u}^{\top} \boldsymbol{C}_{xx} \boldsymbol{u} + \boldsymbol{\lambda} \left( 1 - \boldsymbol{u}^{\top} \boldsymbol{C}_{yy} \boldsymbol{u} \right).$$
(6)

The KKT conditions assert that for the optimal  $(u^*; \lambda^*)$ , it holds that  $C_{xx}u^* = \lambda^* C_{yy}u^*$ , which is a generalized eigenvalue problem. Equivalently, one can rewrite

$$\boldsymbol{C}_{yy}^{-1}\boldsymbol{C}_{xx}\boldsymbol{u}^* = \lambda^*\boldsymbol{u}^* \tag{7}$$

suggesting that  $u^*$  is an eigenvector of  $C_{yy}^{-1}C_{xx}$  associated with eigenvalue  $\lambda^*$ . Respecting the constraint  $(u^*)^{\top}C_{yy}u^* = 1$ , the objective (5a) reduces to

$$(\boldsymbol{u}^*)^{\top} \boldsymbol{C}_{xx} \boldsymbol{u}^* = \lambda^* (\boldsymbol{u}^*)^{\top} \boldsymbol{C}_{yy} \boldsymbol{u}^* = \lambda^*.$$
 (8)

It is clear now that the optimal objective value of problem (5) is equal to the largest eigenvalue of  $C_{yy}^{-1}C_{xx}$ , and the optimal solution  $u^*$  is the corresponding eigenvector.

For d > 1, one finds the d (right) eigenvectors of  $C_{yy}^{-1}C_{xx}$  that correspond to the d largest eigenvalues as the first d dPCs. For ease of implementation, the proposed dPCA approach for contrastive data exploration is summarized in Alg. 1. Concerning dPCA, a couple of remarks are in order.

Algorithm 1 Discriminative principal component analysis.

- 1: **Input:** Nonzero-mean target and background data  $\{\hat{x}_i\}_{1 \le i \le m}, \{\hat{y}_j\}_{1 \le j \le n};$  number of dPCs *d*.
- Remove the mean from {\$\u03c8 x\_i\$} and {\$\u03c9 y\_j\$} to yield centered data {\$x\_i\$}, and {\$y\_i\$}.
- 3: Construct the sample covariance matrices:

$$oldsymbol{C}_{xx} \coloneqq rac{1}{m}\sum_{i=1}^m oldsymbol{x}_ioldsymbol{x}_i^ op, ext{ and } oldsymbol{C}_{yy} \coloneqq rac{1}{n}\sum_{j=1}^noldsymbol{y}_joldsymbol{y}_j^ op$$

- 4: **Perform** SVD on matrix  $C_{yy}^{-1}C_{xx}$ .
- 5: **Output** the *d* (right) singular vectors corresponding to the *d* largest singular values.

**Remark 1.** When there is no background data, with  $C_{yy} = I_D$ , dPCA boils down to PCA. On the other hand, when there are multiple background datasets, one can first combine them into a single one, and then apply dPCA. Other twists will be discussed in the journal version of this paper.

**Remark 2.** Performing dPCA on  $\{x_i\}$  and  $\{y_j\}$  can be seen as performing PCA on the transformed data  $\{C_{yy}^{-\top/2}x_i\}$  to yield  $v^*$ , followed by re-transformation  $u^* = C_{yy}^{-1/2}v^*$ . The new data can be understood as the data obtained after removing the "background" component vectors from the target data.

**Remark 3.** Inexpensive power or Lanczos iterations [14] can be employed to compute the principal eigenvectors in (4).

## 3.2. dPCA vis-à-vis cPCA

Consider again the constrained form of dPCA (5) and its Lagrangian (6). Using Lagrange duality, when choosing  $\alpha = \lambda^*$ in (2), cPCA maximizing  $u^*(C_{xx} - \lambda^*C_{yy})u$  is equivalent to  $\max_{u \in \mathbb{R}^D} \mathcal{L}(u; \lambda^*) = u^\top (C_{xx} - \lambda^*C_{yy})u + \lambda^*$ , which is exactly dPCA. In other words, cPCA and dPCA are equivalent when  $\alpha$  in cPCA is carefully set as the optimal dual variable  $\lambda^*$  for the constrained form (5) of dPCA, namely the largest eigenvalue of  $C_{yy}^{-1}C_{xx}$ .

It will be useful for further analysis to focus on simultaneously diagonalizable matrices  $C_{xx}$  and  $C_{yy}$ , that is

$$C_{xx} := U^{\top} \Sigma_x U$$
, and  $C_{yy} := U^{\top} \Sigma_y U$  (9)

where  $U \in \mathbb{R}^{D \times D}$  is unitary and simultaneously decomposes  $C_{xx}$  and  $C_{yy}$ , while diagonal matrices  $\Sigma_x$ ,  $\Sigma_y \succ 0$  hold accordingly eigenvalues  $\{\lambda_x^i\}_{1 \leq i \leq D}$  of  $C_{xx}$  and  $\{\lambda_y^i\}_{1 \leq i \leq D}$  of  $C_{yy}$  on their main diagonals. It clearly holds that  $C_{yy}^{-1}C_{xx} = U^{\top}\Sigma_y^{-1}\Sigma_x U = U^{\top} \text{diag}(\{\frac{\lambda_x^i}{\lambda_y^i}\}_{1 \leq i \leq D})U$ . Looking for the first *d* dPCs boils down to taking the *d* columns of *U* associated with the *d* largest eigenvalue ratios among  $\{\frac{\lambda_x^i}{\lambda_y^i}\}_{1 \leq i \leq D}$ .

On the other hand, the solution of cPCA under a given  $\alpha$ , or the first d cPCs of  $C_{xx} - \alpha C_{yy} = U^{\top} (\Sigma_x - \alpha \Sigma_y) U = U^{\top} \text{diag} (\{\lambda_x^i - \alpha \lambda_y^i\}_{1 \le i \le D}) U$ , are found as

the *d* columns of U that correspond to the *d* largest numbers in  $\{\lambda_x^i - \alpha \lambda_y^i\}_{1 \le i \le D}$ . In the ensuing section, we show that when given data obey certain models, dPCA is LS optimal.

## 4. OPTIMALITY OF dPCA

Adopting a bilinear (factor analysis) model, PCA describes the (non-centered) data  $\{\mathring{y}_j\}_{1 \le j \le n}$  as

$$\mathring{\boldsymbol{y}}_j = \boldsymbol{m}_y + \boldsymbol{U}_y \boldsymbol{\psi}_j + \boldsymbol{e}_{y,j}, \quad 1 \le j \le n$$
(10)

where  $m_y$  is a location vector,  $U_y \in \mathbb{R}^{D \times D}$  has orthonormal columns;  $\{\psi_j\}_{1 \le j \le n}$  are the coefficients, and  $\{e_{y,j}\}_{1 \le j \le n}$  zero-mean random variables. The unknowns  $m_y$ ,  $U_y$ , and  $\{\psi_i\}_{1 \le j \le n}$  can be estimated using the LS criterion as [15]

$$\min_{\substack{\boldsymbol{m}_{y}, \{\boldsymbol{\psi}_{j}\}\\ \boldsymbol{U}_{y}^{\top}\boldsymbol{U}_{y}=I}} \sum_{j=1}^{n} \left\| \mathring{\boldsymbol{y}}_{j} - \boldsymbol{m}_{y} - \boldsymbol{U}_{y} \boldsymbol{\psi}_{j} \right\|_{2}^{2}.$$
 (11)

whose solution is given as [15, 10]:  $\boldsymbol{m}_y^* := (1/n) \sum_{j=1}^n \hat{\boldsymbol{y}}_j$ ,  $\boldsymbol{\psi}_j^* := (\boldsymbol{U}_y^*)^\top (\hat{\boldsymbol{y}}_j - \hat{\boldsymbol{m}}_y), \forall 1 \leq j \leq n, \text{ and } \boldsymbol{U}_y^* \text{ stacks up as}$  its columns the eigenvectors of  $\boldsymbol{C}_{yy} := (1/n) \sum_{j=1}^n \boldsymbol{y}_j \boldsymbol{y}_j^\top$ , to form  $\boldsymbol{C}_{yy} := \boldsymbol{U}_y^* \boldsymbol{\Sigma}_y (\boldsymbol{U}_y^*)^\top$ , where  $\boldsymbol{y}_j := \hat{\boldsymbol{y}}_j - \boldsymbol{m}_y^*$  is the centered data. For notational brevity, the superscript \* shall be dropped when clear from the context. Wlog, let  $\boldsymbol{U}_y := [\boldsymbol{U}_b \ \boldsymbol{U}_n]$  be partitioned such that  $\boldsymbol{U}_b \in \mathbb{R}^{D \times k}$  corresponds to the first k PCs of  $\boldsymbol{C}_{yy}$ , which capture most background component vectors.

In the context of discriminative data analysis, we *assume* that the target data share some PCs with the background data (say  $U_b$  of (10)), and has additionally a few (say d) PCs that capture patterns specific to the target data but are less significant than the PCs in  $U_b$ . For simplicity, consider d = 1, and model  $\{\hat{x}_i\}$  as

$$\mathring{\boldsymbol{x}}_{i} = \boldsymbol{m}_{x} + [\boldsymbol{U}_{b} \ \boldsymbol{u}_{s}] \begin{bmatrix} \boldsymbol{\chi}_{b,i} \\ \boldsymbol{\chi}_{s,i} \end{bmatrix} + \boldsymbol{e}_{x,i}, \quad 1 \leq i \leq m \quad (12)$$

where  $m_x$  is the mean of  $\{ \hat{x}_i \}_{1 \le i \le m}$ ; and assuming  $k + d \le D$ ,  $U_x := [U_b \ u_s] \in \mathbb{R}^{D \times (k+1)}$  has orthonormal columns, where  $U_b$  describes the component vectors present both in the background as well as target data, while  $u_s \in \mathbb{R}^{D \times 1}$  captures the patterns of interest that are present *only* in the target data. Our goal is to obtain this discriminative subspace  $U_s$  given solely the two datasets. By modeling this distinctly informative component  $\chi_{s,i}u_s$  in (12) explicitly as an outlier vector, it is also possible to employ robust PCA which boils down to solving a nonconvex optimization problem [10].

Likewise, remove the mean  $\mathbf{m}_x := (1/n) \sum_{i=1}^n \mathring{x}_i$  from the target data yielding  $\mathbf{x}_i := \mathring{x}_i - \mathbf{m}_x$ . Recalling  $C_{yy}^{1/2} := \sum_y^{1/2} U_y$ , consider the transformed data model for  $1 \le i \le m$ :

$$C_{yy}^{\top/2} \boldsymbol{x}_{i} = C_{yy}^{\top/2} \begin{bmatrix} \boldsymbol{U}_{b} \ \boldsymbol{u}_{s} \end{bmatrix} \begin{bmatrix} \boldsymbol{\chi}_{b,i} \\ \boldsymbol{\chi}_{s,i} \end{bmatrix} + C_{yy}^{\top/2} \boldsymbol{e}_{x,i}$$
$$= \chi_{s,i} C_{yy}^{\top/2} \boldsymbol{u}_{s} + C_{yy}^{\top/2} \boldsymbol{e}_{x,i} := \chi_{s,i} \tilde{\boldsymbol{u}}_{s} + \tilde{\boldsymbol{e}}_{x,i} \quad (13)$$



Fig. 1. dPCA versus PCA on semi-synthetic data.

where  $C_{yy}^{\top/2} U_b$  vanishes due to the orthogonality of columns of  $U_y^* = [U_b U_n]$ ,  $\tilde{u}_s := (U_y^*)^\top u_s$ , and  $\tilde{e}_{x,i}$  is a zero-mean random variable. Similar to (11), the LS optimal estimate  $\tilde{u}_s^*$ is given by the first principal eigenvector of

$$\tilde{\boldsymbol{C}}_{xx} := (1/m) \sum_{i=1}^m \boldsymbol{C}_{yy}^{\top/2} \boldsymbol{x}_i \boldsymbol{x}_i^\top \boldsymbol{C}_{yy}^{\top/2} = \boldsymbol{C}_{yy}^{\top/2} \boldsymbol{C}_{xx} \boldsymbol{C}_{yy}^{1/2}.$$

Hence, the discriminative PCs can be recovered from  $\tilde{u}_s$  as  $u_s^* := C_{yy}^{1/2} \tilde{u}_s^*$ , which coincides with solutions of problem (3) or (4), and establishes the LS optimality of dPCA.

## 5. NUMERICAL TESTS

In this section, the performance of dPCA is assessed relative to PCA and cPCA [4] on a synthetic and a real dataset. In the first experiment, (semi-)synthetic target and background data were generated from real images. Specifically, the target data were constructed using 2,000 handwritten digits 6 and 9 (1,000 for each) of size  $28 \times 28$  from the MNIST dataset [16] superimposed with 2,000 frog images from the CIFAR-10 dataset [17]. The raw  $32 \times 32$  frog images were converted to grayscale and randomly cropped to  $28 \times 28$ . The background data were built with 3,000 resized images only, which were sampled randomly from the remaining frog images.

We performed PCA on the target data only. The results of the target images embedded onto the first two PCs and two dPCs are depicted in the left and right panels of Fig. 1, respectively. Clearly, PCA is unable to discover the two digit subgroups. This is because the obtained two PCs are likely associated with the background component vectors within the target images, namely features depicting frog images rather than handwritten digits. On the contrary, two clusters emerged in the plot of dPCA, demonstrating its efficacy over PCA in discriminating unique features of one dataset from the other.

The capability of dPCA in discovering subgroups is further tested on real protein expression data. In the second experiment, the target data consist of 267 points, each recording 77 protein expression measurements for a mouse suffering Down Syndrome [18]. There were 135 mice with drugmemantine treatment as well as 134 without treatment. The control data comprise such measurements from 135 healthy mice. The 135 control mice are likely to exhibit similar natural variations (due to e.g., sex and age) as the target mice, but without the differences that result from Down Syndrome. For



Fig. 2. Discovering subgroups in mice protein expression data.

cPCA, the designed spectral-clustering algorithm was implemented for selecting 4 from a list of 15 logarithmically spaced values between 0.001 and 1,000 [4]. The simulated results are presented in Fig. 2, with red circles denoting mice with treatment and blue diamonds the other mice. PCA reveals that the two types of mice follow a similar distribution in the space spanned by the first two PCs; see the left bottom plot in Fig. 2. The separation between the two groups of mice becomes clear when dPCA is applied. At the price of runtime (15 times more than dPCA), cPCA with properly *learnt* parameters ( $\alpha = 3.5938$  and 27.8256) can work well too.

## 6. CONCLUSIONS

This paper put forth a novel approach termed dPCA for discriminative analytics, namely for discovering the most informative features that are specific to one dataset but are also distinct from some other correlated datasets. The resultant algorithm amounts to solving a generalized eigenvalue problem. Comparing to existing alternatives, dPCA bypasses parameter tuning, and it incurs complexity required to perform only one SVD. It is provably optimal in the LS sense provided that the background component vectors are present in the target data. Simulated tests using (semi)-synthetic images and real protein expression data corroborated the merits of the developed approach. Investigating dPCA using kernels and over graphs constitutes meaningful future research directions.

### 7. REFERENCES

- H. Abdi and L. J. Williams, "Principal component analysis," *Wiley Interdiscip. Rev. Comput. Stat.*, vol. 2, no. 4, pp. 433–459, Jun. 2010.
- [2] F. Karl Pearson, "LIII. On lines and planes of closest fit to systems of points in space," *Phil. Mag.*, vol. 2, no. 11, pp. 559–572, 1901.
- [3] H. Hotelling, "Analysis of a complex of statistical variables into principal components," *J. Educ. Psychol.*, vol. 24, no. 6, pp. 417–441, Oct. 1933.
- [4] A. Abid, V. K. Bagaria, M. J. Zhang, and J. Zou, "Contrastive principal component analysis," arXiv:1709.06716, 2017.
- [5] S. Garte, "The role of ethnicity in cancer susceptibility gene polymorphisms: The example of CYP1A1," *Carcinogenesis*, vol. 19, no. 8, pp. 1329–1332, Aug. 1998.
- [6] J. B. Kruskal, "Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis," *Psychometrika*, vol. 29, no. 1, pp. 1–27, Mar. 1964.
- [7] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, Dec. 2000.
- [8] B. Scholkopf, A. Smola, and K. B. Muller, *Kernel principal component analysis*. Berlin, Heidelberg: Springer Berlin Heidelberg, 1997, pp. 583–588.
- [9] I. D. Schizas and G. B. Giannakis, "Covariance eigenvector sparsity for compression and denoising," *IEEE Trans. Signal Process.*, vol. 60, pp. 2408–2421, May 2012.
- [10] G. Mateos and G. B. Giannakis, "Robust PCA as bilinear decomposition with outlier-sparsity regularization," *IEEE Trans. Signal Process.*, vol. 60, no. 10, pp. 5176– 5190, Oct. 2012.
- [11] J. Chen and I. D. Schizas, "Online distributed sparsityaware canonical correlation analysis," *IEEE Trans. Signal Process.*, vol. 64, no. 3, pp. 688–703, Feb. 2016.
- [12] —, "Distributed information-based clustering of heterogeneous sensor data," *Signal Process.*, vol. 126, pp. 35–51, Sep. 2016.
- [13] R. A. Fisher, "The use of multiple measurements in taxonomic problems," *Ann. Eugenics*, vol. 7, no. 2, pp. 179–188, Sep. 1936.
- [14] Y. Saad, Iterative Methods for Sparse Linear Systems. SIAM, 2003.

- [15] B. Yang, "Projection approximation subspace tracking," *IEEE Trans. Signal Process.*, vol. 43, no. 1, pp. 95–107, Jan. 1995.
- [16] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, Nov. 1998.
- [17] A. Krizhevsky, "Learning multiple layers of features from tiny images," in *Master's thesis*, Department of Computer Science, University of Toronto, 2009.
- [18] C. Higuera, K. J. Gardiner, and K. J. Cios, "Selforganizing feature maps identify proteins critical to learning in a mouse model of down syndrome," *PloS ONE*, vol. 10, no. 6, p. e0129126, Jun. 2015.