# An Evaluation of Supervised Dimensionality Reduction For Large Scale Data

**Nancy Jan Sliper**

Advanced Chemical Metallurgy, Norwegian University of Science and Technology, Trondheim, Norway.
nancy.jan@ntnu.no, nancyjan@aol.com

**Abstract** – Experimenters today frequently quantify millions or even billions of characteristics (measurements) each sample to address critical biological issues, in the hopes that machine learning tools would be able to make correct data-driven judgments. An efficient analysis requires a low-dimensional representation that preserves the differentiating features in data whose size and complexity are orders of magnitude apart (e.g., if a certain ailment is present in the person's body). While there are several systems that can handle millions of variables and yet have strong empirical and conceptual guarantees, there are few that can be clearly understood. This research presents an evaluation of supervised dimensionality reduction for large scale data. We provide a methodology for expanding Principal Component Analysis (PCA) by including category moment estimations in low-dimensional projections. Linear Optimum Low-Rank (LOLR) projection, the cheapest variant, includes the class-conditional means. We show that LOLR projections and its extensions enhance representations of data for future classifications while retaining computing flexibility and reliability using both experimental and simulated data benchmark. When it comes to accuracy, LOLR prediction outperforms other modular linear dimension reduction methods that require much longer computation times on conventional computers. LOLR uses more than 150 million attributes in brain image processing datasets, and many genome sequencing datasets have more than half a million attributes.

**Keywords** – Linear Optimum Low-Rank (LOLR), Linear Discriminant Analysis (LDA), Canonical Correlations Analyses (CCA), Principal Component Analysis (PCA), Partial Least Squares (PLS).

## I.  INTRODUCTION

The science and technology of predicting statistical correlations using labeled learning data, known as supervised learning, has allowed a broad range of fundamental and applied discoveries, from detecting indicators in omic data to identifying objects from images. Classification is a kind of supervised learning in which a classifier estimates the "classes" of a new input (for instance, by forecasting sex from Magnetic resonance imaging scanning). Fisher's Linear Discriminant Analysis (LDA) is one of the most fundamental and basic techniques to classification. For a classifier, LDA offers a number of extremely desired qualities. Firstly, it is founded on straightforward geometric rationale: when the input is Gaussian, the averages and deviations include all of the information, hence the best classifier employs both of them. Secondly, LDA may be used to solve issues with several classes. Thirdly, theorems ensure that, underneath the Gaussian assumption, When the small subset $n$ is huge and $p$ dimensionality is small, LDA [1] is the best classifier for the problem. Finally, the methods used to put it into action are quite effective.

On the other hand, modern scientific datasets address classification challenges that did not exist in Fisher's Day. The complexity of databases, in particular, is expanding at an alarming rate. A full genome or connectome, for illustration, might include hundreds and thousands of characteristics or dimensions in today's raw data. However, the larger samples have not increased in lockstep. Because many traditional statistical procedures were created with a "p greater than N" condition in mind, this "P greater than n" dilemma is a non-starter for them. When p is greater or equals to "$n$", executing LDA is like fitting a line to a center: there are arbitrarily many absolutely brilliant fits (all line that go through the juncture), but there is no way of knowing which one is "better." A lack of extra limits will result in a system that will overfit, rather than reducing noise rather than classifying it as such. We also seek methods that are able to handle the data's intricacy, are immune to anomalies, and are rapid to calculate. To overcome these "p is greater than n" issues, many complimentary solutions have been tried.

PCA is the first and maybe most extensively utilized approach. PCA has been cited over 4,000 occasions in 2018, as per PubMed. Other approaches, such as machine learning, regression trees, and sparse learning, that earned 500, 1200, and 2000 clicks, respectively, attracted much less attention in the media. As a result, PCA seems to be the most common high-dimensional issue solver. By reducing the dimensions of the dataset to only those components with the largest variance, PCA "pre-processes" the database PCA is a totally unsupervised dimension reduction strategy, which implies that it does not utilize category label whereas learning low-dimensional frameworks, hence amounting in a lower level of classifier performance. Non-linear manifold learning methods expand PCA, however they frequently don't incorporate class-labeled data, and they

don't do well on large datasets. Even though deep learning, such as using (supervised) auto-encoders, is the most current type of non-linear manifolds learning, these techniques are still poorly understood, involve several parameters to setup, and sometimes do not provide interpretable outputs. In addition, deep learning is vulnerable to the "big data" problem, which happens whenever the sample numbers are more than the complexity of the issue.

There are two methods used to standardize or punish supervised approaches, such as institutionalized LDA and Canonical Correlations Analyses (CCA) [2]. The problem is that in the p > n case, such strategies may be much overfit, lack theoretical foundation in such factors, and have many knobs to transforms; and those that are costlier computationally. Furthermore, PLS (Partial Least Square) [3] is a prevalent method in this group, despite the fact that it does not have solid theoretical certainties or a strong construction. The third most prevalent option for overcoming the "curse of dimensionality" is to use sparse methods. However, precise answers are technologically impossible, and exact solution can only be guaranteed theoretically under extremely strong constraints, and those presumptions are highly brittle. As a result, a gap exists: no known technique can categorize multi-class big datasets with thousands of variables while achieving powerful conceptual assurances, favourable and easily understandable practical performances, and a flexible, durable, and sustainable implementations.

To resolve this problem, we created XOX, an approach for introducing class-conditional period estimations, using LOLR as the most basic examples. The essential idea underlying LOLR is that we could utilize the averages and deviations from every class together (like in CCA and LDA), but without needing more dimension than observations (like in PCA) or imposing severe sparseness constraints. LOLR generates the best low-dimensional reconstruction compared to CCA, LDA, PCA and basic techniques based on random matrix concepts when the input is a Gaussian. Irrespective of the complexity of the attributes, the quantity of observations, or the number of planes in which we projection, this is valid under reasonably relaxed expectations.

The dominance of technologies inferred using XOX methodology LOLR, (ii) a version of XOX known as Quadratic Optimum QDA (QOQ) that allows flexibility of the category covariance matrices, and (iii) a rigorous derivative of LOLR named rLOLR—over other methodologies mathematically in a diverse array of simulation model configurations, including a few that do not follow the presumptions then demonstrated in a diverse array of modelled configurations. Lastly, we demonstrate that on numerous 500 GB neuroscience dataset as well as several multi-gigabyte genomes data, LOLR obtains better precision at lower dimensionality while just taking a few moments on a single computer. This research presents an evaluation of supervised dimensionality reduction for large scale data. We provide a methodology for expanding Principal Component Analysis (PCA) by including category moment estimations in low-dimensional projections. Linear. This paper has been organized as follows: Section II presents the methodology for this research. Section III presents the results of the research, analysis of the XOX framework, and real data applications and benchmarks. Section IV presents a discussion of the results, which Section V draws conclusions to the research.

## II. METHODOLOGY

*Supervised Manifolds Learning*

Fig. 4 depicts a general technique for supervised manifolds training, which is discussed here.

- Step (a): Collect or choose a high-dimensional data training example. For clarity, we utilize the MNIST data, which is one of the most widely used test datasets. The photos in this collection are of handwritten numerals 0 through 9. Each picture is replicated by a $28 \times 28$ matrices, resulting in a dimensionality of $p = 282 = 784$ for the data. We subsample the dataset to pick $n = 300$ samples of the integers 3, 7, and 8 since we are inspired by the $n \ll p$ situation (100 each).

- Step (b): Learn how to use a "reflection" to reduce the dimensionality of high-dimensional datasets. PCA and other manifold learning approaches achieve this by ignoring which pictures belong to which digits (the "classifiers"), but LDA and minimalist techniques attempt to utilize the labels. LOLR is a supervised linear manifolds training approach that learns projection that are linear combination of the initial sample data using the class labels.

- Step (c): Project high-dimensional information into the learnt lower-dimensional space using the learnt projection. This stage requires the knowledge of a projections that can be performed to fresh (testing) data samples whereby the real class labels are unknown. In most cases, non-linear manifold learning techniques cannot be used in this fashion. LOLR, on the other hand, may project fresh sample in such a manner that the data is divided into categories.

- Step (d): Create a classifier using the data's low-dimensional representation. A good classifier reliably assigns the right label to as numerous objects as feasible. As a consequence of applying LDA to the low-dimensional dataset learned by LOLR, an accurate classifier may be generated.

*The Geometric Intuition of Linear Optimum Low-Rank (LOLR)*

We analyze the basic high-dimensional categorization scenario to develop insight for circumstances when LOLR works well and when it does not. We look at n illustrations $(x_i, y_i)$, where $x_i$ represent the p-dimensional topology vectors while $y_i$ represents the binary class labels (either 0 or 1). We presume that input from both categories is equally probable, that both categories have same identical covariance (all characteristics are stochastic with unity variation), and the only variation

between both the categories is their means. The ideal low-dimensional projection in this case may be calculated mathematically: it is the linear combination of the comparison of means and the reverse covariance matrices, often known as Fisher's LDA. Machine learning techniques may be used to determine the parameters when the data distribution is missing, as it is in all actual data issues. However, since the resolution to the underlying statistical issue is not provided for $n < p$, the predicted correlation coefficient will not be invertible, necessitating a different strategy. PCA is widely used to develop a low-dimensional representations, as previously noted. The combined sampling distribution and the pool observed covariance matrix are used in PCA. The top d eigenvalues of the pooled multivariate normality are used in PCA projections after the pooled averages has been removed (in this way omitting all classifications).

LOLR, on the other hand, employs class-conditioned averages and class-centered correlations. By using Fisher's LDA, this approach is able to outperform PCA in terms of productivity. LOLR is created in the following way for a two-class conundrum:

a) Compute sample averages of every class
b) Approximate the variation between averages
c) Computing the class-based covariance matrix, which is, computing the matrix after a subtraction of the class average for every point
d) Computing the eigenvector of the class-conditionally-based covariance
e) Concatenating the variation of the averages with topmost $d - 1$ eigenvector of the class-based covariance.

Class-oriented covariance matrix [4] is important to notice since it affects sample pooling covariance matrix matrices, which are impacted by the fluctuation of class averages. Moreover, as seen in the methodology Section II, the class-based covariance matrix is the same as the rrLDA (Reduced Rank LDA). Section III presents the results of this research.

### III. RESULTS

*Flexibility and Accuracy of the XOX Framework*

We use simulation to test the adaptability and reliability of XOX in a way that goes past hypothetical promises. We select 100 training datasets, each with 100 variables, for three distinct situations; hence, Fisher's LDA is unable to address the issue (so many methods of overfitting may be found). To reflect the data into a low-dimensional domain, we investigate a variety of approaches, like PCA, rrLDA, PLS, Randomized Projection (RP), and CCA. Following data projection, we possible facilitate the learning of either QDA (the third case) or LDA (the first two cases), which generalise LDA by enabling every category has its own correlation matrix. On held-out data, we analyze the classification error for every scenario.

Fig. 1 illustrates a 2D scatter plot (left) and the mis-categorization rates against dimensionality (right) for each simulation. Afterwards, LOLR stands for a robust approximation of locality (the class median, which corresponds to the centralized moments whenever the population is distributed symmetrically), and the truncated single value decomposes to the estimates of the second moments. Robust locality estimate makes minimal variations whenever robust estimates were not needed, and empirically advances the performances of simulations and actual-data samples whenever robust approximations were acceptable. Other options included directly using robust estimates of both moment two and moment one. We do not utilize robust approximates the $2^{nd}$ time, as robust projection of the $2^{nd}$ time alludes to the standardized numerical package requisite $d < n$ that is not suitable for big datasets. The upper C – 1 of LOLR embedment dimension is in correspondent to its performances after a projection onto category-conditional averages, whereas rrLDA corresponds to its performance after projections onto class-conditional covariance matrices (CCMs). Fig. 1 (a) depicts a three-category generality aspect of the Trunk examples as depicted in Fig. 5 (b).

LOLR, in contrast to ROAD, which is limited to two courses, may be simply extended to accommodate several classes. Typical of modern biomedical datasets, Fig. 1 (b) displays a two-class sample with several outliers. However, LOLR and PLS both perform well and successfully locate embedded dimension. Unlike PCA or rrLDA, Fig. 1 (c) illustrates a case that should be in opposition to LOLR's use of both. For this reason, LOLR incorporates additional noise-inducing parameters. While LOLR posits that the correlations in both classes is same, the correlations in class-conditional correlations are bidirectional. Compared to PCA, LOLR does a respectable job, but not nearly as well.

An alternative QDA classifier known as QOQ also employs XOX, but instead of using LDA to get the eigenvectors, it sorts the singular values of the concatenated eigenvectors before determining the eigenvectors for each class. For categorization purposes, QOQ is marginally more efficient than PCA [5]. While the first few dimension (those covered by the divergence of the means) are unspecific, the succeeding dimensions (those covered by a larger range of values) are more insightful (the class-conditional covariance).

Extensive variations of LOLR In each of the three scenarios, QOQ attains mis-categorization rates, which are more comparable to or lower compared to other techniques. We may see from these findings that simple XOX modifications of LOLR that contain alternative or robust moment estimations can significantly outperform other projection techniques. There are other systems that don't allow for this kind of flexibility, which is a stark difference.
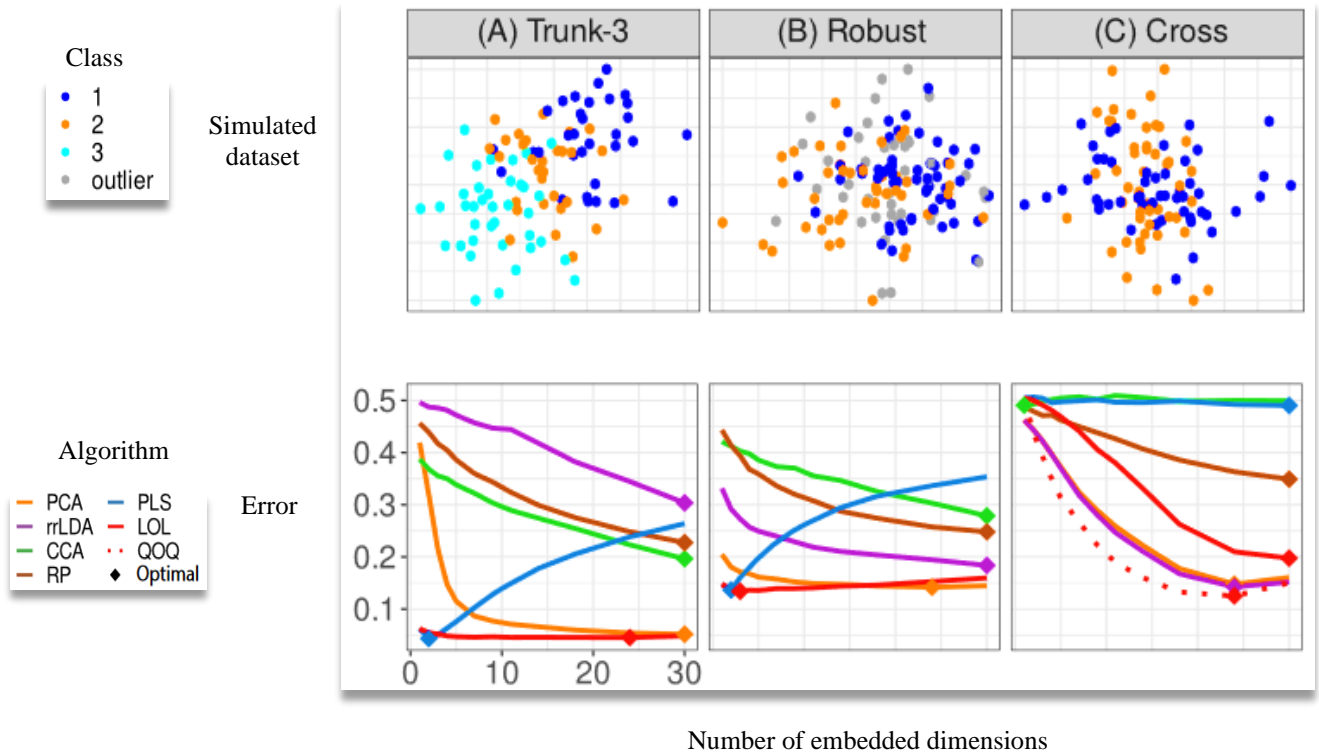
Number of embedded dimensions

**Fig 1.** Three (3) simulation illustrating the accuracy and flexibility of the XOX framework in configurations beyond the present theoretical foundations

Training data has a 100-sample size and 100-dimensionality sample size. Two of the 100 dimensions of the categorization job are shown in the first row of data points. Variations in the bottom-row misclassification rate in relation to the number of projected dimensions [6]. To classify the embedded data, we use LDA and QDA classifiers. The settings for the simulation are:

(a) Trunk-3: This is an alternative to Fig. 5 (b) in which there are three distinct categories.
(b) Robust: The projection matrix is estimated using a large number of outliers in the sample. Laughter is resistant to outliers because of the high confidence in the first moment estimate.
(c) Cross: The means of the two groups are identical, but their covariances are orthogonal. Using the QDA classifier, points are categorized. Because it incorporates the covariance for each class, QOQ outperforms other approaches as one may anticipate, the LOLR variant. For the most part, LOLR, or a generalization of it, is superior than alternative techniques in all situations and dimensions.

*Computational Scalability and Efficiency of XOX*

As dimensionality increases (in the millions or billions), the most common bottleneck is the inability to do anything with the data [7]. It is shown that LOLR has good computing efficiency and scalability in the most basic test case: 1000 samples per class of spherically symmetric Gaussian. This is because, unlike iterative programming approaches now needed for problems of the sparse or dictionary learning kind, LOLR allows for a closed-form solution to be used. As a companion to the R package used in these figures, we produced Flash-LOLR (FlashLOLR), a fast and scalable LOLR implementation that has R bindings.

LOLR's scalability is facilitated by four qualities. In the first place, LOLR is a linear function of sample size and dimension Second, "semi-external memory" (Fig. 2 (a), the red lines (dashed) highlighted that LOLR is in line with multiple cores) makes LOLR easily parallelizable. It's worth mentioning that LOLR doesn't add any significant computational burden to PCA (orange dashed line). A third way that LOLR may improve its performance is by the implementation of randomized approximation techniques for eigen decompositions (Fig. 2 (a), orange line). If you use highly scares RP other than eigenvector, Flash Low-rank Fast Linear (FlashLFL) embeddings, deliver the orders of speeding magnitudes. LOLR is now using nested hyper-parameter selections, which allows all lower-dimensional projections to be instantly available upon estimation of the d-dimensional projections. When a penalty term is fine-tuned, it creates a new optimization issue for each distinct value of a parameter. In other words, LOLR's computational complexity [8] is a 0 (one zero). $(\frac{npd}{Tc})$, where $n$ represents the sample sizes, $p$ represents the data dimensionality, $d$ represents the projection dimensionality, $T$ represent the threads, and $c$ represents the projection sparsity.

(a)                                                        (b)
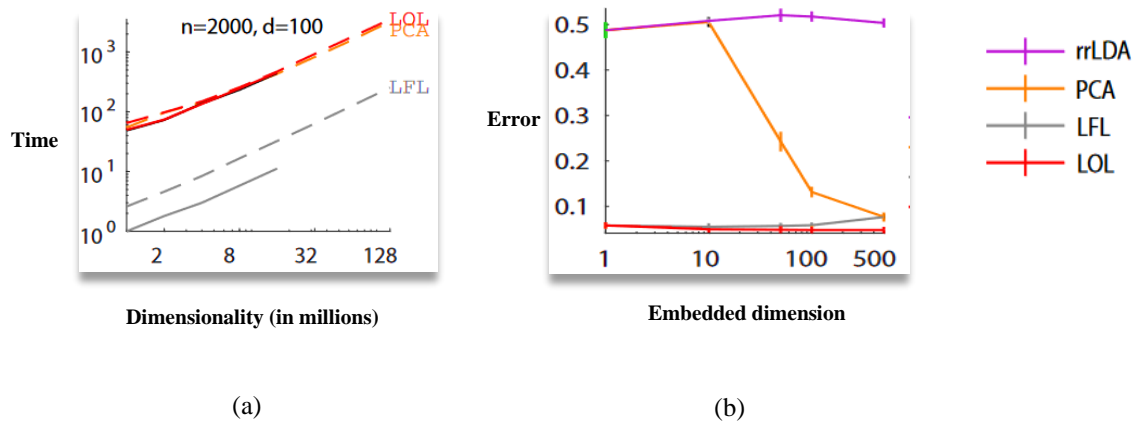
**Fig 2.** Computational scalability and efficiency of LOLR utilizing "n" as 2000 samples from the symmetrical Gaussian datasets (spherical)

(a) A 500 GB dataset requires 46 minutes to identify the projection, but LFL takes just 3 minutes (semi-external memory is shown via dashed lines).

(b) Both LFL and LOLR have lower error than PCA and rrLDA in this situation, and this holds true regardless of whether the projection dimensions are computed using a random technique or not. LFL retains the great performance of LOLR despite the randomization of PCA in (B) despite the scalability increases that may be made to PCA in (a).

For PCA and rrLDA, this simulation configuration is excellent since the first principal component is comprised of the mean difference vector. While rrLDA relies on probability and PCA takes 500 parameters to even approximate the reliability that LOLR obtains with only one dimensions, LOLR and LFL both reach near-optimal precision. PCA would also benefit from a randomized method, however it is important to stress that LFL preserves the outstanding results of LOLR in contrast to PCA, with the added bonus of increased computing efficiency.

*Real Data Applications and Benchmarks*

Aside from the aforementioned simulations, real data may give a different viewpoint on the performance attributes of various algorithms than can be gained through simulation studies. One set of issues stems from computed tomography, the other from genomes, as we explain here. This is a classification issue that we can look at in both circumstances. Alcalde-Barros et al. [9] often use substantiative preprocessing procedures to lower the size of the data in order to categorize participants. Preprocessing methodologies continue to be debated, hence there is no de facto "standard" for best practices in preprocessing data. Parametric modeling and downsampling are common in traditional techniques, which often entail a long processing chain with several phases. Thus, we study the potential of immediately categorizing on the practically raw, high-dimensional dataset.

Individual-specific identifies for individual's sex were allotted to $n > 800$ participants from 5 distinct processing locations in the Consortium for Reliability and Reproducibility (CRR). There are more than 150 million variables in each brain volume, and each dataset contains between 42 and 400 samples (60 GB to 600 GB of data).

A large genomics dataset, comprising of 340 people, is also taken into account, consisting of 144 individuals with nonmetastatic cancer and 196 control subjects, of whom 198 men and 142 women are females and males respectively. Samples from cancer patients are examined for aneuploidy (abnormal chromosomal numbers) by aligning them to more than 750,000 amplicons spread throughout the genome. No additional preprocessing is required for the amplicon counts. We're interested in either sex or age categorization, therefore we have two options.

Once the embedding matrices have been computed for each of the issues above, we use LDA to categorize low-dimensional approximations of the training data. It is therefore forecasted and categorized based on the trained classifier and embedded matrix, and average cross-validated errors are identified over the various dataset folds. Every strategy's optimal multiplicity is defined by the number of embedded parameters with the lowest mean cross-validation error for each challenge. Using Cohen's Kappa, we are able to compare the performance of different classification strategies by standardizing their effectiveness between zero and one (a random chance classifier may be compared to this classifier) (the classifier completes flawlessly). Lastly, for every forecast method, we subtract (PCA) from (PCA) to compute the effect size (embed).

FlashLOLR is the only algorithm that can execute on this dataset with a single thread on a normal desktop machine. Only LOLR outperforms PCA on all tasks in Fig. 3 (a). Each dataset's mean cross-validated misclassification frequency for the LDA classifier is shown in Fig. 3 (b), after the datasets were projected using the stated embedding approach. There are no difficulties that LOLR can't handle. The Wilcoxon signed-rank statistic shows that LOLR outperforms all other methods (all p values are 0.008). In all datasets, LOLR achieves an average misclassification rate of between 5 percent and 15 percent, which is in line with the performance we and others get using substantial data processing and downsampling on comparable

datasets. By not preparing at all, Chapman et al. [10] may bypass the discussed preprocessing difficulties by just employing LOLR to the data set in its original dimension.
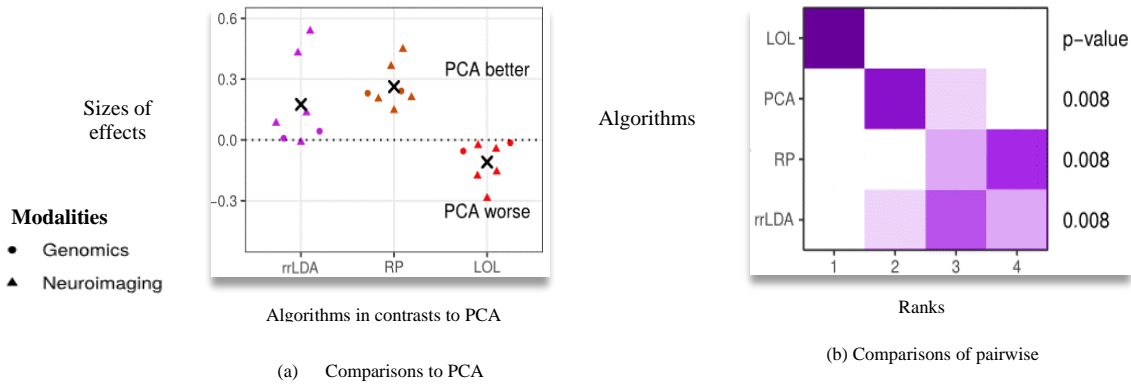


**Fig 3.** Comparisons of several dimensionality reduction techniques into two actual datasets (genomics and neuroimaging)

(a) Comparisons to PCA: At an optimum number of embedded metrics, the numbers of these embedded components with least classifier, each approach's classifiers are shown in a series of plots. LDA, PCA and LDA embed are two embedding approaches that are often used in conjunction with one another, and the effect size is one way to compare the two. The black "x" denotes the sample size-weighted mean impact of PCA compared to the other approach on a particular dataset. It's always better to use LOLR than PCA or any other method.

(b) Comparisons of pairwise: For every dataset, frequency histogram illustrate the respective ranks for every embedding technique, with 1 representing the best relative classification accuracy and 4 denoting the lowest relative classification efficiency, after embedded using the methodology described above. A Wilcoxon signed-rank test shows that LOLR gives a considerable advantage over rival techniques on all benchmark issues (n = 7).

## IV.  DISCUSSION

With the use of class-conditional instances, we've shown a reasonably simple method for increasing the efficacy of supervised learning algorithms with large datasets. (Large data collection having at least as many dimensions as a randomly selected sample). LOLR, in particular, makes advantage of both the discrepancy in averages and the class-centered covariances, allowing it to beat PCA and other supervised linear classification systems in a broad range of cases while incurring no extra processing expense. Using robust estimation technique and/or class-specific covariance parameter estimation, straightforward extensions allow robust and non-linear variants. Our open-source approach grows up to terabyte datasets with ease. Furthermore, the idea may be used to hypothesis testing as well as regression.

PLS and CCA are two techniques that are widely used in these situations. When p n, CCA is identical to rrLDA, which is irrelevant in this case. CCA and rrLDA are not equal when p n; nonetheless, CCA demonstrates the "maximal data stacking issue" in such cases. Basically, each class's points are all projected onto the same spot. This causes substantial overfitting of datasets, amount to poorer empirical findings in various cases we looked at (CCA's initial metric is often worse than the comparison of averages). In spite of these inadequacies, PLS lacks robust conceptual certainties and a simple geometric understanding. When there are outliers or the invariant border is quadratic, CCA and PLS do not allow for straightforward adaptations, unlike XOX (see Fig. 1). Furthermore, XOX consistently beats both of these techniques in all scenarios, sometimes significantly. On a wide range of simulations, XOX beats CCA, for example. No scalability or parallelized applications of these ideas can be found readily on the Internet (see Fig. 2). To address these additional optimization issues, one might apply gradient descent with sanctions, but this would require tuning the penalty value, which would be operationally expensive. PLS or CCA might be conducted efficiently on the enormous neuroimaging databases or the amplicon-level genomics dataset using generally available technology.

Several earlier studies have tackled similar issues. Fisher's LDA was the first to be used with PCA in the famous Fisherfaces article (as an alternative to PCA in this paper). The authors demonstrated the value of projections data utilizing PCA before categorizing with LDA via a series of numerical tests. This work is expanded by including a supervised element in the first projection. We also present the geometric understanding for why and when adding supervision is beneficial, with multiple instances illustrating its superiority as well as conceptual assurances formulating when LOLR exceeds PCA. Similar ideas may be found in the "sufficient complexity reductions" research, which uses a different approach that necessitates the complexity to be less than the random sample. Theoretical restrictions on linear and quadratic classification performances have recently been generated using communication-inspired classifiers; however, they do not rigorously compare various projections, and our constraints are more general and tighter. The aforementioned strategies, on the other hand, are not able to handle billions of qualities at once. Recent big data software is built to handle billions upon billions of data. In biological

sciences, however, tens or even hundreds of data and thousands or millions of characteristics are significantly more typical (e.g., genomics or connectomics).

While most diverse learning approaches have great conceptual and practical performances, they are often completely unsupervised. When it comes to classifying information, they tend to use a one-dimensional approach instead of paying attention to the labels. This method may be problematic if the discriminating dimension and the maximum variance orientations in the trained manifolds are not allocated (Fig. 4 for distinct samples). As for supervised learning, the non-linear manifold learning techniques learn mappings from the source data to environments that are of low-dimension but not projections showcasing that novel dataset cannot be projected with ease onto environments of low-dimensions, which is necessary. However, although deep learning methods may be easily supervised, they sometimes need enormous data, lack conceptual certainties, or are unreachable "black box", which are unsuitable for various therapeutic application. As a result, supervised scaled image compression approaches with powerful conceptual assurances for classification techniques constraints developed for large datasets are few. Even though random forests address many of these concerns, no implementation currently exists that can handle thousands of variables, and retrieved features typically perform no more compared to PCA on big data.
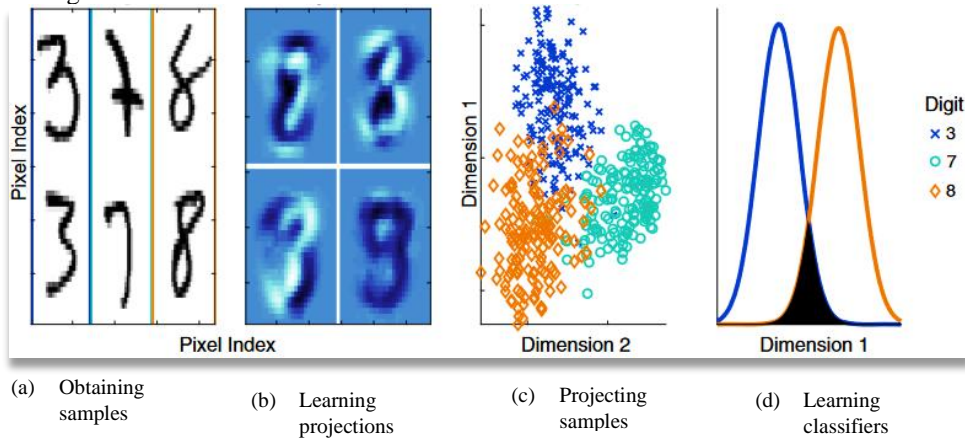


| (a) Obtaining samples | (b) Learning projections | (c) Projecting samples | (d) Learning classifiers |

**Fig 4.** Schematic representation of LOLR as a supervised learning approach

(a) Collect or choose a high-dimensional data training example: More than 784-dimensional images of the numbers 3, 7, and 8 from the MNIST collection are used as training examples (The use of boundary colors is only for the sake of visualisation).

(b) Learn how to use a "reflection" to reduce the dimensionality of high-dimensional datasets: LOLR taught the first four projections matrices. Each image is a composite of the preceding examples.

(c) Project high-dimensional information into the learnt lower-dimensional space using the learnt projection: Fresh (test) samples of 500 were taken and color-coded based on the two most learned dimensions (A). Data from three separate clusters that were LOLR-projected.

(d) Create a classifier using the data's low-dimensional representation: Learning a classifier from low-dimensional data. 3 and 8 of the test samples' estimated distributions show that they can be readily separated by linear techniques following LOLR projections (after data were projected into two dimensions and classified using LDA) (The digit is indicated by the color of the line). An algorithm's purpose is to reduce the region that is filled in by the anticipated error rate. In this high-dimensional actual data scenario, LOLR performs well.

According to Fig. 5, LOL attains near-optimum performance for the 3 multi-variate Gaussian distribution, each having 100 samples in 100 distinct dimensionalities. To categorize 10,000 additional samples, we project onto the top three variables for each technique, and then utilize LDA to do so. From top to bottom, here are the six rows of data: A scatter plot of the measurement points in the first 2 dimensions, with 1 and 0 signified by blue and orange dashes, is indicated in the first row. Class 1 and 0 are illustrated in dotted and solid lines correspondingly. In order to determine the amount of the mistake, we need to know how much overlap there is between the distributions. The estimated criterion for each approach is shown by the vertical black line. A number of approaches are available, including reduced rank LDA (rrLDA), and PCA that projects into eigenvectors to the test category-conditional covariances; ROAD that is developed for this framework; and Bayes optimum classifiers.

(a) Stacked cigars: The mean variation vectors are connected to the vector of maximum variances and concentrated in one dimension, hence making it appealing for PCA, sparse techniques and rrLDA. Based on this context, the findings are similar in distinct approaches, and basically optimum.

(b) Trunk: The average variation vectors are impertinent as per the vectors of maximum discrepancy; The performances of PCA are bas while rrLDA are viewed as the sparse methodologies. LOL could possibly recover the real magnitudes, accomplishing approximately optimum performances.

(c) Rotating trunks (similar to (b), but datasets are rotational, in this instance, LOL's performance is notable. It should be noted that LOL is nearer to Bayes optimum as visualized in three cases.
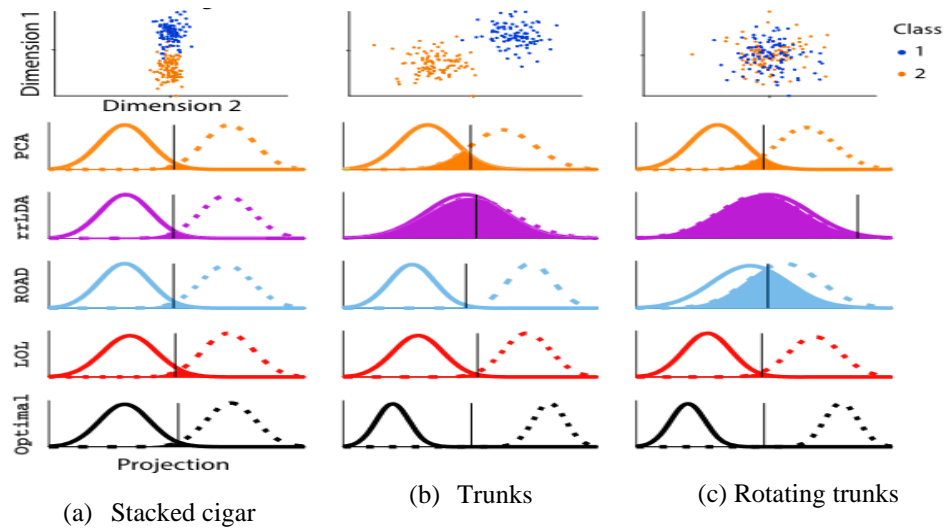


(a) Stacked cigar     (b) Trunks     (c) Rotating trunks

**Fig 5.** LOL's close-optimal performances for 3 multi-variate Gaussian distribution, having 100 examples each in 100 distinct dimensionalities

While projection chasing and empirical risk reduction both attempt to solve the optimization problem, they go about it in a completely different approach. It is important to note that these techniques are confined since they are potentially disposed to localized minima and need time-consuming recursive iterative processes. Before the execution of PCA on the remnant features, techniques such as stronger criticism thresholding and feature selection processes such as higher criticism denoising may be used to effectively filter dimensions. These strategies might be used in combination with LOLR in ultra-high-dimensional environments. In addition, another recently suggested supervised PCA variant employs Hilbert-Schmidt independence necessities for learning embedment. While this linear projection may be used with the difference in means, our theory predicts that this combination will lead to better results in a broad variety of situations. This work might be enhanced by estimating a Gaussian mixture framework for every class, instead of a typical Gaussian for every class on the subspaces that span in accordance to wide-range Gaussian models.

## V. CONCLUSION

An important aspect of unsupervised manifold learning may be transformed into supervised manifold learning by adding class-conditional moment estimations. The first three, LOLR, QOQ, and rLOLR, demonstrated XOX's capacity to adapt to both theoretical and benchmark contexts. Higher-order moments, kernel techniques, ensemble methods like multiscale and random forests methodologies are also of urgent interests since they incorporate additional non-linearities. Even in the Gaussian model with extra (very moderate) restrictions, LOLR outperforms or is on equal footing with the rrLDA. PCA is also outperformed or on par with this product (Other linear projection methods are available as well.). We may apply this to all possible dimensionalities and number of dimensions in our projections given enough big data sizes. Conclusively, these requirements are projected to expand directly to the dimensional increase. Detailed arguments and theorems supporting the aforesaid conclusion are offered under "Methodology" in Section II. PCA's covariance matrix contains more classification information than LOLR's covariance matrix does, hence PCA is the preferable choice when comparing the two. Mean difference vectors should be included in projection matrixes rather than ignored, assuming the same assumptions that drive PCA, as stated above. An upgrade or worse performance does not correlate with the array of options in an extracted features or how many samples there are, but rather the number of projections dimension $d$.

## References

[1]. C. S. Pun and M. Z. Hadimaja, "A self-calibrated direct approach to precision matrix estimation and linear discriminant analysis in high dimensions," Comput. Stat. Data Anal., vol. 155, no. 107105, p. 107105, 2021.

[2]. L. Cope, D. Q. Naiman, and G. Parmigiani, "Integrative correlation: Properties and relation to canonical correlations," J. Multivar. Anal., vol. 123, pp. 270–280, 2014.

[3]. Z. Cai and B. Chen, "Least-squares method for the Oseen equation: Least Squares For Oseen's Problem," Numer. Methods Partial Differ. Equ., vol. 32, no. 4, pp. 1289–1303, 2016.

[4]. M. L. Stein, J. Chen, and M. Anitescu, "Difference Filter Preconditioning for Large Covariance Matrices," SIAM J. Matrix Anal. Appl., vol. 33, no. 1, pp. 52–72, 2012.

[5]. M. Pal, N. K. Mandal, and M. L. Aggarwal, "A-optimal designs for optimum mixture in an additive quadratic mixture model," Statistics (Ber.), vol. 51, no. 2, pp. 265–276, 2017.

[6]. J. Busa and I. Polaka, "Variability of classification results in data with high dimensionality and small sample size," Inf. Technol. Manag. Sci., vol. 24, pp. 45–52, 2021.

[7]. S. Mainali et al., "An Information-theoretic approach to dimensionality reduction in data science," Int. J. Data Sci. Anal., 2021.

[8]. N. Fischer and C. Ikenmeyer, "The computational complexity of plethysm coefficients," Comput. Complex., vol. 29, no. 2, 2020.

[9]. A. Alcalde-Barros, D. García-Gil, S. García, and F. Herrera, "DPASF: a flink library for streaming data preprocessing," Big Data Anal., vol. 4, no. 1, 2019.

[10]. A. Chapman, P. Missier, G. Simonelli, and R. Torlone, "Capturing and querying fine-grained provenance of preprocessing pipelines in data science," Proceedings VLDB Endowment, vol. 14, no. 4, pp. 507–520, 2020.