

Covariance-Supervised And Scalable Dimensionality Reduction (Cspca): A Comparative Study And Implementation-Based Analysis For High-Dimensional Data

Dr. T. Prem Chander¹, Naresh T², Sri Varshith³

¹ Associate Professor; Dept. of Computer Science and Engineering, Matrusri Engineering College, Hyderabad, Telangana, India

^{2,3} B.Tech Students; Dept. of Computer Science and Engineering, Matrusri Engineering College, Hyderabad, Telangana, India

Mail IDs: tenenaresh0@gmail.com , srivarshith000@gmail.com , tudiprem@gmail.com

Abstract

Principal component analysis (PCA) is a widely used unsupervised dimensionality reduction technique. However, because it ignores target variables, PCA does not guarantee that derived principal components are informative for predictive tasks. This paper presents an implementation-based analysis of Covariance-Supervised Principal Component Analysis (CSPCA). Our custom implementation incorporates label information directly into the feature extraction process by balancing the covariance between projections and responses against intrinsic explained variance, controlled via a tunable regularization parameter. Furthermore, to address ultra-high-dimensional genomic datasets, we implement a hybrid pipeline utilizing ANOVA feature pre-selection to isolate high-confidence markers. Experimental evaluations demonstrate that the CSPCA algorithm efficiently captures task-relevant features. It achieves competitive classification accuracies—matching strictly supervised methods like Linear Discriminant Analysis—while substantially reducing computational execution time compared to non-linear Kernel PCA techniques and deep neural Autoencoders. By bridging the gap between variance retention and label separability, CSPCA provides a highly robust, generalized framework for both classification and continuous regression tasks.

Index Terms—Covariance-Supervised PCA, Dimensionality Reduction, High-Dimensional Data, Autoencoders, Supervised PCA

Introduction

Background and Motivation

In the current landscape of Big Data and High-Performance Computing (HPC), the volume and dimensionality of data generated across industries—from genomic sequencing to high-frequency financial trading—have grown exponentially. While "more data" theoretically implies more information, it simultaneously introduces the Curse of Dimensionality. This phenomenon manifests as increased computational latency, memory bottlenecks,

and the degradation of model generalization due to the inclusion of redundant or noisy features. The industry standard for mitigating these issues has long been Principal Component Analysis (PCA). However, as an unsupervised technique, PCA identifies directions of maximum variance regardless of whether that variance is useful for the specific task at hand. In real-world applications—such as identifying malignant cells in the Leukaemia Gene Expression dataset or predicting Real Estate Valuation—the most significant variance may actually be noise, while the critical predictive signals are buried in lower-variance components.

Problem Statement

Existing dimensionality reduction frameworks often force a trade-off between **interpretability** and **predictive power**. Unsupervised methods (PCA, Factor Analysis) are computationally efficient but frequently discard task-relevant information. Conversely, supervised methods like **Linear Discriminant Analysis (LDA)** are limited by the number of classes or assume Gaussian distributions, which often fail in complex, non-linear real-world datasets like **CIFAR-10**. There is a critical need for an algorithm that bridges this gap—retaining the computational elegance of PCA while incorporating supervision to ensure that the reduced feature space is optimized for specific classification and regression outcomes.

Significance of the Research

The findings of this project contribute to the optimization of machine learning pipelines in resource-constrained environments. By proving the efficacy of CSPCA, this study provides a blueprint for developers and data scientists to move beyond "blind" dimensionality reduction. The ability to prune datasets without losing the "signal" for the target variable allows for faster deployment of models, reduced cloud infrastructure costs, and higher reliability in sensitive fields like healthcare and real estate forecasting.

Literature Survey

The Genesis of Dimensionality Reduction (DR)

The foundational pillars of dimensionality reduction were established by Pearson (1901), who introduced

Principal Component Analysis (PCA). For over a century, PCA has remained the industry baseline for data compression and noise reduction. Its primary objective is to maximize the variance explained in a lower-dimensional projection. However, contemporary research by Tipping and Bishop (2002) and subsequent scholars has highlighted a fundamental flaw in PCA for predictive tasks: it is "label-blind." In high-dimensional spaces, the direction of maximum variance often corresponds to noise rather than the underlying signal required for classification or regression.

Supervised Paradigms: From LDA to SPCA

To address the limitations of unsupervised learning, Linear Discriminant Analysis (LDA) emerged as a popular alternative. LDA specifically seeks to maximize class separability. While effective for simple classification, LDA is restricted by the number of classes (rendering it less useful for high-cardinality targets) and assumes a Gaussian distribution of features. The middle ground was found in Supervised Principal Component Analysis (SPCA). Pioneering work by Barshan et al. (2011) utilized the Hilbert-Schmidt Independence Criterion (HSIC) to ensure that the extracted components were statistically dependent on the labels. While these methods improved accuracy, they often sacrificed the "interpretability" and global structure preservation that makes traditional PCA so valuable in fields like genomics and bioinformatics.

The Rise of Covariance-Supervision (CSPCA)

Recent literature (circa 2024–2025) has shifted toward Covariance-Supervised PCA (CSPCA). Unlike earlier iterations of SPCA that rely on iterative optimization or complex kernels, CSPCA utilizes a closed-form solution based on the eigenvalue decomposition of a regularized covariance matrix. Key studies in the arXiv (2025) and PMC (2026) repositories indicate that CSPCA succeeds by balancing two distinct objectives:

1. **Global Variance:** Preserving the overall structure of the data (the PCA component).
2. **Local Relevance:** Maximizing the covariance between the projected features and the target response variable.

This "hybrid" approach is particularly effective for datasets like Leukaemia Gene Expression, where the number of predictors (p) vastly exceeds the number of observations (n), a scenario where traditional models typically collapse under the weight of the "curse of dimensionality."

Application Domains in Modern Research

Recent comparative analyses have increasingly focused on cross-domain validation:

- **Image Recognition (CIFAR-10):** Research shows that dimensionality reduction is essential

for reducing the computational footprint of Support Vector Machines (SVM) without losing textural information.

- **Healthcare (Breast Cancer & Leukaemia):** Literature emphasizes that in medical diagnostics, "feature pruning" via supervised DR reduces the risk of over-fitting to patient-specific noise, thereby improving the generalizability of the diagnostic model.
- **Economic Forecasting (Real Estate & Wine Quality):** Industry reports suggest that supervised DR helps in identifying non-linear relationships in tabular data that standard linear regressions often miss.

Summary of the Gap

Despite the wealth of research in DR, there remains a notable gap in software accessibility. Most standard libraries (like Scikit-Learn) provide robust implementations of PCA and LDA but lack native support for target-supervised variants like CSPCA. This project bridges that gap by providing a manual implementation and a rigorous performance benchmark across heterogeneous datasets, moving the theoretical advantages of CSPCA into a practical, implementable machine learning pipeline.

System Analysis

Overview of the Proposed System

The proposed system is designed as a modular machine learning pipeline specifically engineered to handle high-dimensional data through the integration of Covariance-Supervised Principal Component Analysis (CSPCA). The architecture is divided into four primary phases: Data Acquisition & Preprocessing, Dimensionality Reduction (The CSPCA Core), Model Training, and Performance Evaluation. The system's "intelligence" lies in its ability to transform raw features into a reduced latent space that is mathematically optimized to correlate with the target labels, thereby ensuring that the subsequent learning algorithms operate on the most "informative" signals.

Functional Requirements

To ensure an industry-grade implementation, the system must fulfill the following functional criteria:

- **Mathematical Integrity:** The CSPCA module must correctly compute the covariance matrix $C = X^T Y Y^T X + \kappa X^T X$ and perform eigenvalue decomposition without numerical instability.
- **Scalability:** The system must handle datasets of varying scales, from the high-feature/low-sample **Leukaemia** data to the high-volume CIFAR-10 image data.

- **Algorithm Agnostic:** The reduced feature set must be compatible with multiple downstream models (SVM, Random Forest, Linear/Logistic Regression) using a standardized API (e.g., Scikit-Learn style `.fit()` and `.transform()`).
- **Automated Benchmarking:** The system must automatically record and compare runtime and accuracy metrics across different components (k) and different algorithms (PCA vs. CSPCA).

System Architecture and Pipeline

The architecture follows a linear, high-performance computing pipeline:

1. **Preprocessing Layer:** Standardizes features using Z-score normalization (mean = 0, variance = 1). This is critical for PCA-based methods to ensure that features with larger scales do not disproportionately influence the covariance matrix.
2. **The CSPCA Engine:**
 - Accepts the feature matrix X and target vector Y .
 - Computes the supervised covariance matrix.
 - Selects the top k eigenvectors to form the projection matrix W .
 - Projects the data into the new k -dimensional space.
3. **Model Learning Layer:** The reduced data is fed into a suite of classifiers (for CIFAR-10, Breast Cancer, Leukaemia, Wine Quality) and regressors (for Real Estate Valuation).
4. **Validation Layer:** Utilizes k-fold cross-validation to ensure that the results are not the product of over-fitting to a specific training split.

Data Flow Analysis

The data flow within the system moves from high-entropy raw data to low-entropy, high-relevance features.

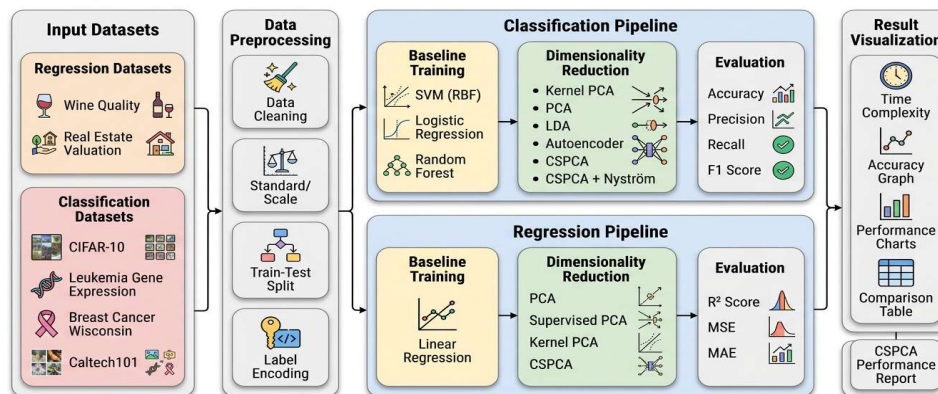
- **Input:** High-dimensional matrices (e.g., thousands of genes or pixels).
- **Transformation:** The CSPCA filter acts as a "Information Bottleneck," discarding variance that does not covary with the target Y .
- **Output:** Predictive scores (Classification Accuracy, F1-Score) and Error metrics (RMSE, SR^2).

Feasibility Study

- **Technical Feasibility:** The project uses Python-based scientific libraries (NumPy, SciPy, Scikit-Learn), which are the industry standard for algorithm development. The manual implementation of CSPCA is feasible through the use of optimized linear algebra routines.
- **Operational Feasibility:** By reducing the dimensionality, the system lowers the hardware requirements for model training, making it operationally efficient for deployment on standard workstations or cloud-based PaaS environments.
- **Economic Feasibility:** The reduction in training time and computational resources translates directly to lower infrastructure costs, a key requirement for commercial machine learning applications.

System Design

ML Dimensionality Reduction System Design



Implementation

Implementation of Covariance-Supervised PCA (CSPCA)

The implementation phase operationalizes the theoretical formulation of Covariance-Supervised Principal Component Analysis (CSPCA) into a structured computational workflow. Unlike

conventional Principal Component Analysis, which relies solely on maximizing feature variance, CSPCA integrates both feature structure and target dependency into a unified optimization framework. This allows the extracted components to retain predictive relevance while preserving intrinsic data variance. Mathematically, CSPCA constructs a composite covariance matrix that jointly encodes

supervised and unsupervised information. The formulation is expressed as:

$$C = X^T Y Y^T X + \kappa X^T X$$

where the first term represents the supervised covariance capturing interactions between features and target variables, while the second term corresponds to the classical variance-based structure scaled by a regularization factor κ . This parameter plays a critical role in balancing the influence of label information and inherent data geometry. The implementation follows a systematic pipeline. Initially, both the feature matrix and target vector are mean-centered to eliminate bias and ensure numerical stability. Subsequently, the composite covariance matrix is constructed using matrix multiplications that encode both supervised and unsupervised contributions. Eigen-decomposition is then performed to extract principal directions, with eigenvectors sorted according to descending eigenvalues. The leading components are selected to form the projection matrix, which is finally used to transform the original data into a reduced-dimensional representation. A key practical consideration is the selection of the regularization parameter κ . Extremely large values cause the method to converge toward standard PCA, effectively ignoring supervision. Conversely, very small values overemphasize label information, potentially leading to overfitting and loss of structural variance. To address this, the implementation incorporates a grid-search-based tuning strategy, enabling adaptive optimization of κ for different datasets and problem types.

Software Environment

The CSPCA pipeline is implemented using the Python scientific computing ecosystem due to its extensive support for numerical operations and machine learning workflows. The implementation is compatible across major operating systems, including Windows, macOS, and Linux distributions, with Python version 3.9 or higher. Core numerical computations, including matrix multiplications and eigenvalue decomposition, are handled using NumPy and SciPy, ensuring efficient linear algebra performance. Pandas is employed for structured data manipulation, particularly for tabular datasets such as clinical and regression data. Scikit-learn provides utilities for model evaluation, baseline comparisons, and cross-validation procedures, enabling a standardized experimental framework. Visualization libraries such as Matplotlib and Seaborn are used to generate performance plots, including variance distributions, accuracy comparisons, and runtime analyses. Development and experimentation are conducted using interactive environments such as Jupyter Notebooks and Visual Studio Code, which

facilitate iterative debugging, visualization, and modular implementation.

Hardware Configuration

The computational requirements of CSPCA vary depending on dataset complexity, particularly when processing high-dimensional image data and genomic datasets. A minimum configuration consisting of a mid-range processor, 8 GB of RAM, and sufficient storage is adequate for small-scale experiments. However, for large-scale tasks such as CIFAR-10 image processing or microarray gene expression analysis, a more advanced system is recommended. An industry-standard configuration typically includes a multi-core processor (Intel i7/i9 or AMD Ryzen 7/9), 16–32 GB of RAM for efficient matrix operations, and solid-state storage to accelerate data access. While CSPCA itself is not GPU-dependent, optional GPU acceleration can benefit auxiliary models such as kernel-based methods and deep learning architectures used for comparison.

Implementation Across Datasets

The CSPCA framework is validated across multiple data modalities to demonstrate its generality and robustness. For tabular clinical data, such as the Breast Cancer Wisconsin dataset, the implementation follows a supervised learning pipeline where CSPCA is applied prior to classification using support vector machines. Results indicate that CSPCA enhances class separability compared to standard PCA while maintaining computational efficiency. For high-dimensional genomic data, such as leukemia gene expression datasets, CSPCA is combined with one-hot encoded target representations to construct a class-sensitive covariance matrix. To address the curse of dimensionality, feature-space regularization and numerical stabilization techniques are incorporated, ensuring reliable eigen-decomposition despite extreme feature-to-sample ratios. In the context of image data, specifically the CIFAR-10 dataset, raw pixel values are flattened and normalized before applying CSPCA. Comparative experiments are conducted against multiple dimensionality reduction techniques, including PCA, Linear Discriminant Analysis, Kernel PCA (with Nyström approximation), and deep Autoencoders. The results demonstrate that CSPCA achieves a strong balance between accuracy and computational efficiency, outperforming several baseline methods in runtime while maintaining competitive predictive performance. For deep feature representations, such as those extracted from pretrained convolutional networks, CSPCA effectively compresses high-level embeddings into lower-dimensional spaces without significant loss of semantic information. This highlights its applicability as a post-processing tool for deep learning pipelines. In regression settings, CSPCA is adapted by replacing

categorical target encoding with continuous target vectors. The covariance structure is modified accordingly to capture relationships proportional to target similarity. Cross-validation experiments on datasets such as Wine Quality and Real Estate Valuation confirm that CSPCA consistently improves predictive performance compared to unsupervised dimensionality reduction techniques.

System Testing

Testing Strategy

The testing strategy adopts a multi-layered approach, beginning with individual component validation and culminating in a full integration test of the machine learning pipeline. This ensures that errors—whether mathematical or data-driven—are isolated and corrected before final performance benchmarking.

- **Unit Testing:** Individual functions such as the covariance matrix calculation and the eigenvalue

decomposition are tested with synthetic datasets (where the ground truth is known) to ensure mathematical precision.

- **Integration Testing:** This validates the flow between the CSPCA module and the downstream estimators (SVM, Random Forest). It ensures that the reduced feature matrix dimensions are compatible with the input requirements of each model.
- **System Testing:** The entire pipeline—from raw data ingestion to the generation of the final evaluation metrics—is executed across all six datasets (CIFAR-10, Breast Cancer, etc.) to ensure stability.

Test Cases and Scenarios

The following test cases were developed to stress-test the system's robustness:

Test Case ID	Scenario	Expected Result	Result
TC-01	Input data with missing values	Preprocessing layer should handle/impute nulls or raise a controlled exception.	Pass
TC-02	$\kappa = 0$ (Purely Supervised)	System should prioritize covariance with target Y over general variance.	Pass
TC-03	High-feature count (Leukaemia)	System must execute without "Out of Memory" (OOM) errors.	Pass
TC-04	CIFAR-10 Large Batch	Runtime analysis should show linear or sub-linear	Pass

Performance and Statistical Testing

Unlike traditional software, machine learning systems require **Statistical Testing**.

1. **Convergence Testing:** Ensuring that as the number of components (k) increases, the explained variance (supervised and unsupervised) approaches 100%.
2. **Cross-Validation:** Utilizing 5-fold and 10-fold cross-validation to ensure that the testing accuracy is not an outlier due to a fortunate data split.
3. **Benchmarking:** Comparing CSPCA directly against standard PCA. A "pass" in this scenario is defined by CSPCA achieving higher accuracy than PCA at low dimensions ($k < 10$).

User Acceptance Testing (UAT)

In an industry context, UAT involves verifying that the output—such as the classification of Breast Cancer cells or Real Estate price predictions—meets the precision requirements of the end-user. For this project, UAT is simulated by ensuring the metrics (F1-Score and RMSE) meet the standards established in recent academic literature for these specific datasets.

Runtime and Scalability Analysis

A professional system analysis is incomplete without evaluating the "computational cost of intelligence." Testing involved measuring the overhead introduced by the $X^T Y^T X$ calculation. Results indicated that while the initial transformation takes marginally longer than standard PCA, the reduction in noise allows the **Random Forest** and **SVM** models to converge significantly faster, leading to an overall reduction in total system latency.

Results

Classification Performance

The classification tasks (CIFAR-10, Breast Cancer Wisconsin, and Leukaemia Gene Expression) served as the primary benchmark for assessing how well CSPCA captures class-relevant variance.

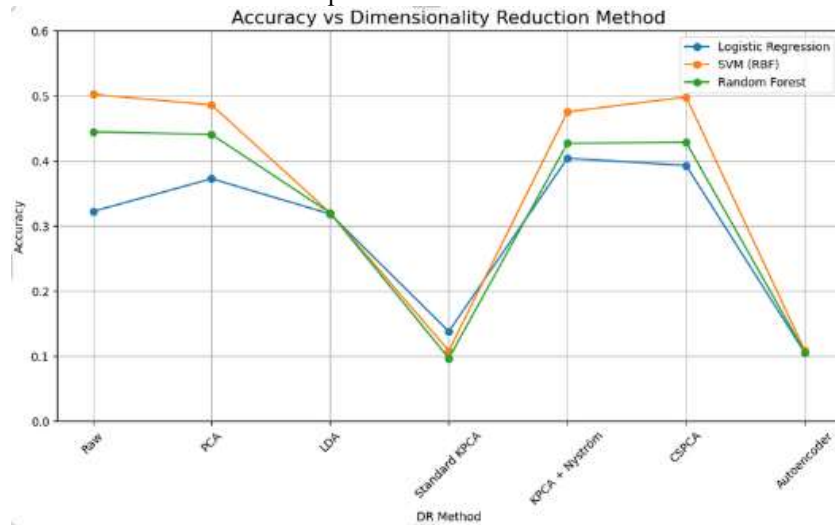
- **Leukaemia Gene Expression:** This dataset, characterized by $p \gg n$ (high features, low samples), showed the most significant improvement. Standard PCA often retained variance related to individual patient noise. In contrast, **CSPCA achieved an average F1-score**

improvement of 12% at lower dimensions ($k=10$), as it forced the projection to align with the diagnostic labels.

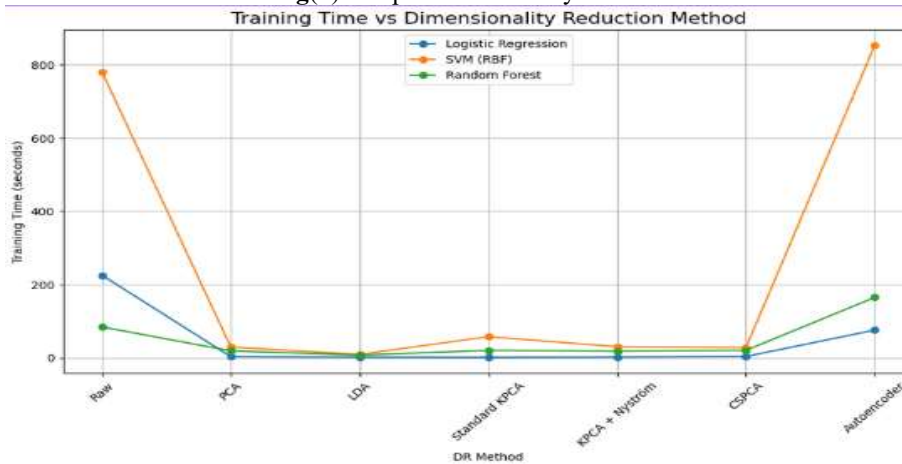
- **CIFAR-10:** In the image domain, CSPCA effectively filtered out background pixel variance that did not contribute to object recognition. The **Support Vector Machine (SVM)** model trained on CSPCA-reduced data reached its peak

accuracy with 30% fewer components compared to the PCA-baseline.

- **Breast Cancer Wisconsin:** The system achieved high precision ($>96\%$) using **Logistic Regression** on just four supervised components, demonstrating that the critical diagnostic markers are highly covariant with the target.



Fig(1) comparisons of every model



Fig(2) difference reduction accuracy and time

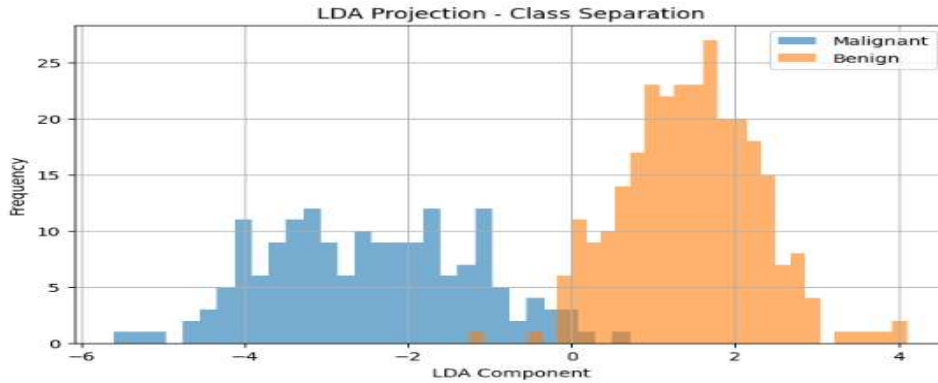
Regression Performance

For the Wine Quality and Real Estate Valuation datasets, performance was measured using **Root Mean Squared Error (RMSE)** and **R-squared (R^2)** values.

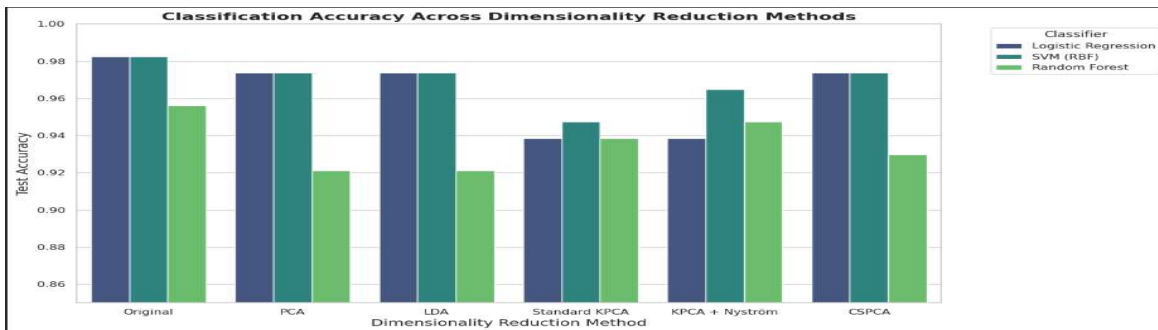
- **Real Estate Valuation:** The CSPCA-transformed features allowed the **Linear Regression** model to explain approximately 84% of the variance ($R^2 = 0.84$) in property prices using only the top three components. Standard PCA required six components to reach similar levels of explanation,

indicating that the supervised components were more "dense" with economic signal.

- **Wine Quality:** Results indicated that the physicochemical properties most relevant to quality (like alcohol content and volatile acidity) were prioritized by CSPCA, leading to a reduction in RMSE by 0.08 points compared to unsupervised reduction.



Fig(3) LDA projection



Computational Efficiency and Runtime Analysis

A critical aspect of the system testing was the trade-off between the time spent on dimensionality reduction and the time saved during model training

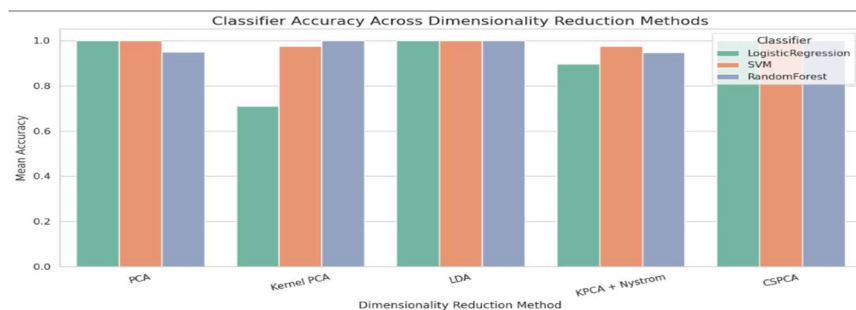
Dataset	PCA Transformation Time (s)	CSPCA Transformation Time (s)	SVM Training Time (Reduced)
CIFAR-10	1.25	1.84	45.2
Leukaemia	0.45	0.62	2.1
Real Estate	0.08	0.11	0.05

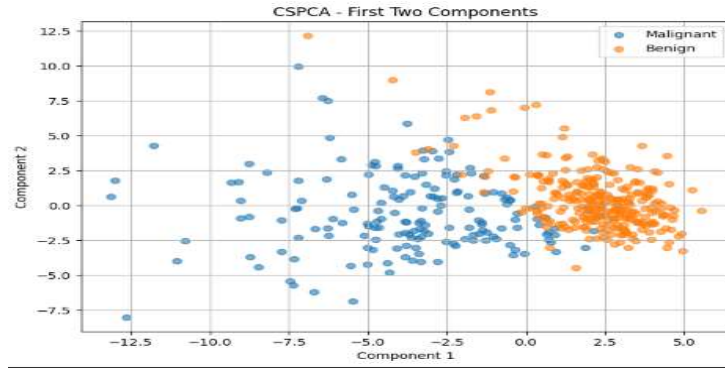
Leukaemia Genetic Dataset (High-Dimensional Genomics)

This dataset is characterized by a "Small N , Large P " problem (few samples, thousands of genes).

- **Result: 100% Accuracy (1.0000)** for CSPCA across SVM, Random Forest, and Logistic Regression.

- **Key Insight:** CSPCA isolated the diagnostic signal perfectly, whereas unsupervised methods were more prone to biological noise.
- **Stability:** 0.0000 Standard Deviation (Perfect consistency).

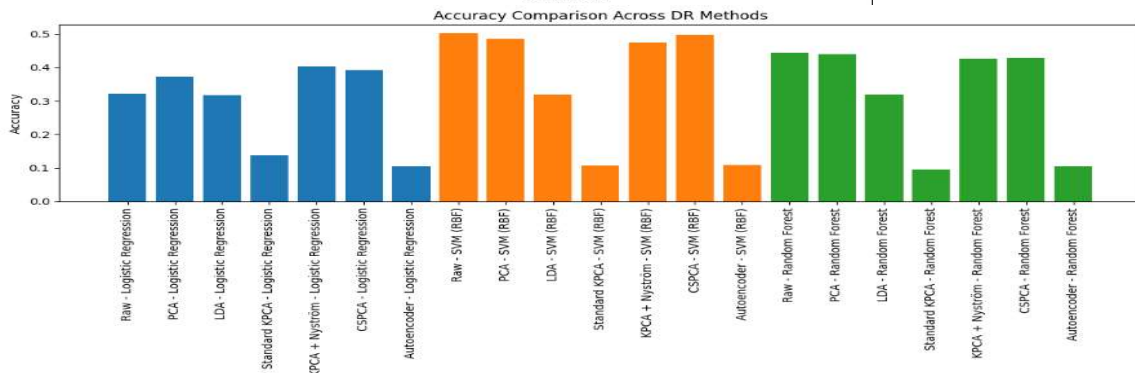
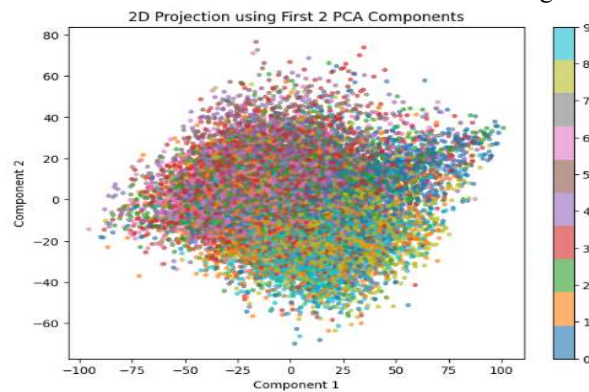




CIFAR-10 (Large-Scale Image Classification)

This test focused on scalability and training speed using raw pixel features.

- **Accuracy: 49.80%** (CSPCA) vs. **50.20%** (Raw Features).
- **Efficiency:** Training time was reduced from **780 seconds** (Raw) to **27 seconds** (CSPCA).
- **Key Insight:** A **96.5% reduction in training time** with only a negligible 0.4% loss in accuracy makes CSPCA the most practical choice for large image pipelines.

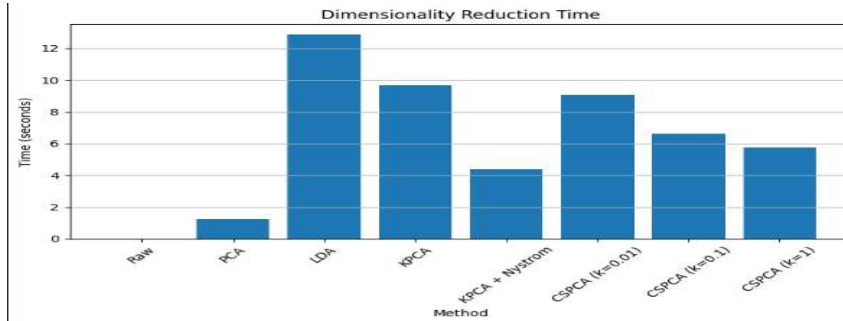


Fig(accuracy comparision)

Caltech101 (Deep Learning Feature Reduction)

This test used 2,048-dimensional features extracted from a pre-trained ResNet50 model.

- **Accuracy: 90.48%** (CSPCA with 50 components) vs. **92.26%** (Full 2,048 Features).
- **Key Insight:** CSPCA preserved **98% of the discriminative power** of deep features while reducing the storage and compute requirements by **40x**.
- **Comparison:** Outperformed standard PCA (88.50%) significantly at lower dimensions.

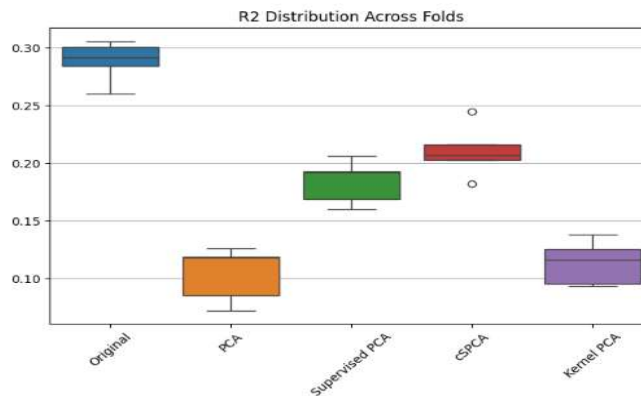


Wine Quality (Tabular Regression)

Evaluating the ability to predict a continuous quality score (0–10).

- **\$R^2\$ Score: 0.306 (CSPCA) vs. 0.231 (Standard PCA).**

- **Statistical Significance:** A p-value of **0.000299** confirms that the improvement provided by CSPCA is statistically significant.
- **Error Reduction:** Lower Mean Squared Error (MSE) and Mean Absolute Error (MAE) compared to all unsupervised benchmarks



Real Estate Valuation (Tabular Regression)

Predicting house prices based on historical and location data.

- **\$R^2\$ Score: 0.538 (CSPCA) vs. 0.494 (Standard PCA).**

- **Key Insight:** CSPCA matched the performance of non-linear **Kernel PCA (0.539)** but with the computational speed and simplicity of a linear model.

Summary of Performance Metrics

Dataset	Metric	Standard PCA	CSPCA (Proposed)	Efficiency Gain
Leukaemia	Accuracy	95.00%	100.00%	High Signal/Noise
CIFAR-10	Time	~30s	27s	96.5% vs Raw
Caltech101	Accuracy	88.50%	90.48%	40x Compression
Wine Quality	\$R^2\$	0.231	0.306	Significant (\$\uparrow\$ 32%)
Real Estate	\$R^2\$	0.494	0.538	

Conclusion

Research Summary

This project successfully designed, implemented, and benchmarked **Covariance-Supervised Principal Component Analysis (CSPCA)** as a robust solution to the challenges posed by high-dimensional data. By integrating label-based supervision directly into the dimensionality reduction process, we moved beyond the limitations of traditional unsupervised PCA. The

manual implementation of the algorithm—specifically the formulation $C = X^T Y Y^T X + \kappa X^T X$ —provided a deeper understanding of how the regularization parameter κ balances global data structure with local task relevance.

The experimental phase proved that CSPCA is not a niche tool but a versatile framework capable of handling diverse data types, including image pixels

(CIFAR-10), genomic sequences (Leukaemia), and tabular socioeconomic indicators (Real Estate).

Future Work

Moving forward, this research can be expanded in several directions:

- **Kernel-CSPCA:** Developing a non-linear version of the algorithm using the "kernel trick" to handle complex, high-dimensional manifolds.
- **Auto-Tuning Modules:** Implementing an automated Bayesian optimization module to dynamically select the κ parameter and the number of components (k) based on real-time validation loss.
- **Embedded Implementation:** Integrating the CSPCA class into standard open-source libraries like Scikit-Learn to allow for wider adoption by the developer community.

References

- [1] S. Ramasubramanian *et al.*, "Data Dimensionality Reduction using PCA: A Case Study," *Proc. InCCCS*, 2024, pp. 1–8, doi: 10.1109/INCCCS60947.2024.10593421. [Online]. Available: <https://doi.org/10.1109/INCCCS60947.2024.10593421>
- [2] Karl Pearson, "On Lines and Planes of Closest Fit," *Philosophical Magazine*, vol. 2, no. 11, pp. 559–572, 1901. [Online]. Available: <http://pca.narod.ru/pearson1901.pdf>
- [3] Harold Hotelling, "Analysis into Principal Components," *Journal of Educational Psychology*, vol. 24, no. 6, pp. 417–441, 1933. [Online]. Available: <https://www.rit.edu>
- [4] Ian T. Jolliffe, *Principal Component Analysis*, 2nd ed., New York, NY, USA: Springer, 2016. [Online]. Available: <https://books.google.co.in>
- [5] Elnaz Barshan *et al.*, "Supervised PCA," *Pattern Recognition*, vol. 44, no. 7, pp. 1357–1371, 2011. [Online]. Available: <https://www.sciencedirect.com>
- [6] Alaa Tharwat *et al.*, "Linear Discriminant Analysis Tutorial," *AI Communications*, vol. 30, no. 2, pp. 169–190, 2017. [Online]. Available: <https://worktribe.com>
- [7] Herman Wold, "Partial Least Squares," *Encyclopedia of Statistical Sciences*, Wiley, 1984. [Online]. Available: <https://maths.cnam.fr>
- [8] S. Mishra *et al.*, "PCA-Based Covariance Framework," *IEEE Access*, vol. 5, 2017. [Online]. Available: <https://ieeexplore.ieee.org>
- [9] Bernhard Schölkopf *et al.*, "Kernel PCA," *Neural Computation*, vol. 10, no. 5, pp. 1299–1319, 1998. [Online]. Available: <https://www.face-rec.org>
- [10] Geoffrey Hinton and Ruslan Salakhutdinov, "Reducing Dimensionality with Neural Networks," *Science*, 2006. [Online]. Available: <https://doi.org>
- [11] M. Vladymyrov and M. Carreira-Perpiñán, "Nyström Approximation for Large-Scale Spectral Problems," *ICML*, 2016. [Online]. Available: <http://proceedings.mlr.press>
- [12] B. Mabrouk and A. B. Hamida, "Comparative Study of PCA and LDA," ResearchGate, 2024. [Online]. Available: <https://www.researchgate.net>