

Dimensionality Reduction

Doluwamu Taiwo Kuye¹

Contents

1	Motivation	2
2	Dimensionality reduction	3
2.1	Application	4
2.2	Challenges	4
3	Evaluation Of Principal Component Analysis method	5
3.1	Iris Dataset	5
3.2	Code	5
3.3	RESULT ANALYSIS	6
4	Conclusion	7
5	Declaration of Originality	8

Abstract: The objectives, approaches, applications, and challenges connected with dimensionality reduction in machine learning are discussed in this article. Overfitting, sparsity, and processing complexity are all issues with high-dimensional datasets. Dimensionality reduction seeks to improve the performance of machine learning models, obtain insights into data structure, and handle sparsity difficulties. Dimensionality reduction is accomplished by the use of several approaches, such as PCA, SVD, NMF, and t-SNE. Dimensionality reduction has applications in data visualization, feature selection, image and speech recognition, and anomaly detection. However, there are difficulties in selecting a strategy, estimating probable information loss, and finding a suitable reduced dimensionality.

¹ Doluwamu-taiwo.kuye@stud.hshl.de

1 Motivation

High-dimensional datasets are common in many scientific and engineering domains, such as computer vision, bioinformatics, physics, and finance. However, the curse of dimensionality, or the inherent difficulty of dealing with high-dimensional data, can cause a variety of issues such as overfitting, sparsity, and computational complexity [Ha09]. As a result, in recent years, dimensionality reduction has been an active study area, with numerous advanced algorithms being developed to meet these issues. One of the primary goals for dimensionality reduction is to improve machine learning model performance by reducing noise and minimizing overfitting [RM17]. In high-dimensional datasets, the number of features can be significantly greater than the number of samples, which can lead to overfitting and poor generalization performance. High-dimensional datasets can be difficult to visualize and analyze; by reducing the dimensionality, we can gain a better understanding of the correlations between the features and the samples. For example, one of the most prominent dimensionality reduction approaches, principal component analysis (PCA), can highlight the most important lines of variation in the data, which can then be used to find patterns and anomalies in the dataset [AW10]. Other techniques, such as t-SNE, can be used to visualize high-dimensional data in low-dimensional space, which can help to identify clusters and subgroups in the dataset [LH08]. In addition, dimensionality reduction can also help to address the issue of sparsity in high-dimensional datasets [G  22]. In Figure 1 we see a showcase of dimensionality reduction.

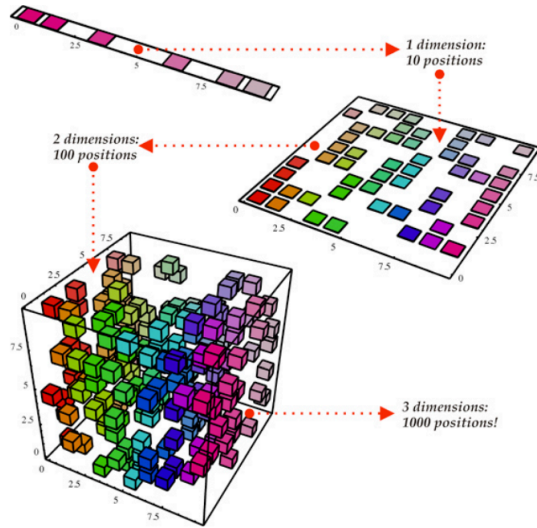


Fig. 1: Dimensionality reduction VISUAL AID [Pi23].

2 Dimensionality reduction

Dimensionality reduction is a machine learning technique that seeks to minimize the number of features or dimensions in a dataset while retaining as much useful information as feasible [G 22]. This technique is useful in various fields such as image and speech recognition, bioinformatics, and finance, where high-dimensional data are common. The importance of dimensionality reduction is rooted in the curse of dimensionality [G 22]. As the number of dimensions increases, so does the number of samples required to cover the space, making it more difficult to discern between samples and perform statistical analysis [Ha09]. Dimensionality reduction helps alleviate this issue by removing redundant or noisy features, increasing computational efficiency, and lowering the danger of overfitting. Dimensionality reduction can be approached in two ways: feature selection and feature extraction. Feature selection is concerned with identifying and removing irrelevant or redundant features based on some criterion, such as correlation or mutual information, whereas feature extraction is concerned with transforming the original features into a lower-dimensional representation that maintains the most relevant information [RM17]. One of the most prominent feature extraction techniques is Principal Component Analysis (PCA). It projects the data onto a lower-dimensional space using linear transformations while maximizing the variance of the projected data [AW10]. Other popular techniques include Singular Value Decomposition (SVD), Non-negative Matrix Factorization (NMF), and t-SNE [LS99]. Dimensionality reduction is a powerful machine-learning method, but it has limitations. One of the most difficult difficulties is identifying the ideal number of dimensions to reduce [BN06]. Dimensionality reduction helps to mitigate the curse of dimensionality while also improving computational performance. In Figure 2 we see another example of dimensionality reduction.

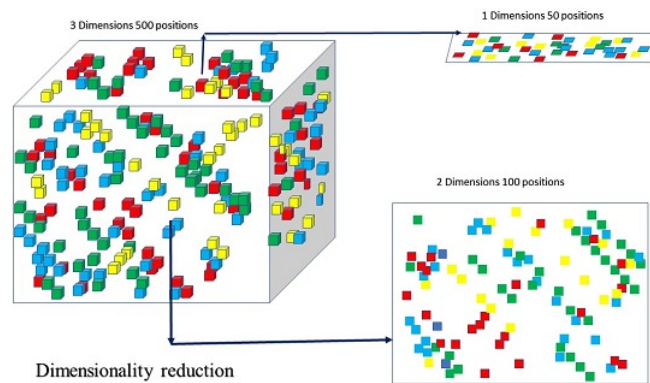


Fig. 2: Depicts dimensionality reduction in 1 and 2 dimensions [R23].

2.1 Application

- Data visualization is one of the most common applications of dimensionality reduction. High-dimensional data can be difficult to visualize and understand, especially when the relationships between its features are complicated. Dimensionality reduction techniques such as PCA and t-SNE can assist in projecting data into a lower-dimensional space while retaining as much structure and relationships as feasible [LH08]. The resulting visualization can assist in the identification of patterns, clusters, and outliers that were not visible in the high-dimensional space. This application has been used in a variety of fields, including biology, finance, and social network analysis [LH08].
- Feature selection is another application of dimensionality reduction. Not all features in the high-dimensional data are useful or informative for the task at hand in some situations. Dimensionality reduction techniques, such as PCA and LASSO, can aid in the identification and removal of redundant or irrelevant features, simplifying the data and boosting the performance of future machine learning models [RM17].

2.2 Challenges

- Technique selection: One of the most difficult aspects of dimensionality reduction is determining what technique is most appropriate for a given dataset. There are numerous dimensionality reduction techniques available, each with its own set of advantages and disadvantages, and no single methodology is universally applicable [Va09].
- Another challenge with dimensionality reduction is the possible loss of information that can occur when a dataset's dimensionality is reduced. While dimensionality reduction can be effective for reducing complex data, it can also result in the loss of vital information that the reduced representation does not capture [RM17].
- Determining the appropriate dimension: Choosing the optimal reduced dimensionality is a difficult task. Several methods have been presented, including sequential testing, bootstrap processes, BIC type criteria, and sparse eigen-decomposition, however, each method has limitations and may not consistently give accurate findings [MZ13].

Despite these challenges and contradictions, dimensionality reduction is an important and commonly utilized technique in a variety of fields. Researchers and practitioners can continue to improve the effectiveness and utility of this strong methodology by carefully considering the strengths and limits of various strategies and developing new methods that address the challenges and contradictions of dimensionality reduction.

3 Evaluation Of Principal Component Analysis method

3.1 Iris Dataset

To evaluate dimensionality reduction, I used the PCA technique to reduce the collection of data from the Iris dataset. The Iris dataset contains measurements of sepal length, sepal width, petal length, and petal width from three different species of iris flowers (Setosa, Versicolor, and Virginica). It is frequently used in machine learning for applications including as classification, clustering, and dimensionality reduction.

	A	B	C	D	E
1	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
2	5.1	3.5	1.4		0.2 setosa
3	4.9	3	1.4		0.2 setosa
4	4.7	3.2	1.3		0.2 setosa
5	4.6	3.1	1.5		0.2 setosa
6	5	3.6	1.4		0.2 setosa
7	5.4	3.9	1.7		0.4 setosa
8	4.6	3.4	1.4		0.3 setosa
9	5	3.4	1.5		0.2 setosa
10	4.4	2.9	1.4		0.2 setosa
11	4.9	3.1	1.5		0.1 setosa
12	5.4	3.7	1.5		0.2 setosa
13	4.8	3.4	1.6		0.2 setosa
14	4.8	3	1.4		0.1 setosa
15	4.3	3	1.1		0.1 setosa
16	5.8	4	1.2		0.2 setosa
17	5.7	4.4	1.5		0.4 setosa
18	5.4	3.9	1.3		0.4 setosa
19	5.1	3.5	1.4		0.3 setosa
20	5.7	3.8	1.7		0.3 setosa

Fig. 3: A snippet of the iris dataset used full data set can be found [here](#).

3.2 Code

The data is in four dimensions, which is difficult for the human brain to visualize; however, we can reduce the dimensionality to two dimensions using the PCA technique. This procedure was carried out utilizing the SCIKIT learning environment and the necessary libraries. Importing the libraries, standardizing the feature values (essential for more accurate results), instantiating our PCA, and visualizing the results are the main steps. This is all shown in my code below in 4.

```
import numpy as np
import matplotlib.pyplot as plt

from sklearn.datasets import load_iris
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

# Load the iris dataset
iris = load_iris()

# Standardize the data
scaler = StandardScaler()
X = scaler.fit_transform(iris.data)

# Apply PCA to reduce the dimensionality of the data to 2 dimensions
pca = PCA(n_components=2)
X_reduced = pca.fit_transform(X)

# Get unique target labels and corresponding names
target_labels = np.unique(iris.target)
target_names = iris.target_names[target_labels]

# Plot the reduced data with labeled points
for target_label, target_name in zip(target_labels, target_names):
    plt.scatter(X_reduced[iris.target == target_label, 0],
                X_reduced[iris.target == target_label, 1],
                label=target_name)

plt.xlabel('PC1')
plt.ylabel('PC2')

plt.legend()

plt.show()
```

Fig. 4: Main code can be found [here](#).

After applying PCA on the Iris dataset and reducing the dimensionality to two dimensions, the resulting scatter plot will show the dataset instances on the PC1 and PC2 axes. PC1 and PC2 are the first and second principal components obtained from PCA, respectively. The principle components are linear combinations of the original features that capture the most significant patterns and variations in the data.

3.3 RESULT ANALYSIS

The plot in 5 demonstrates how the original high-dimensional data (four features: sepal length, sepal width, petal length, and petal width) was translated into a lower-dimensional space (two primary components: PC1 and PC2). The axes of the figure correspond to these key components. The scatter plot depicts the projections of the original data onto the dataset's two most significant directions of variation, PC1, and PC2. For example, we can observe on the plot that setosa varies significantly from versicolor and virginica, indicating that setosa has different values and patterns than the other two. Verisicolor is also discernible from virginica by a small margin. This variation can be attributed to the fact that PCA seeks the directions (principal components) that explain the greatest amount of variance in the data, and these directions can reveal intrinsic differences across classes.

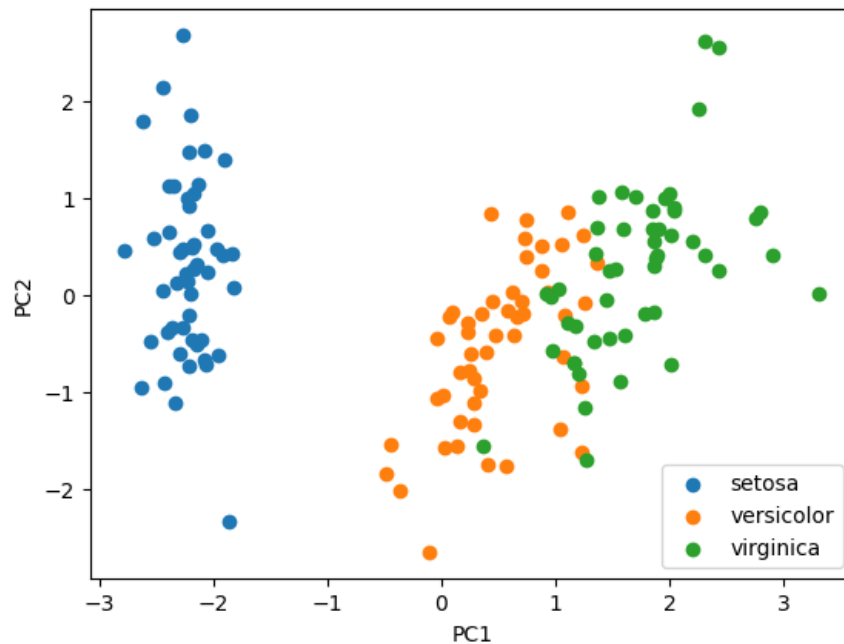


Fig. 5: PCA PLOT

My simulation's explained variance ratio is [0.72962445, 0.22850762]. It is essential to note that the explained variance is a key factor for determining how much of the variance in data is retained after reduction. The percentage of my PC1 and PC2 is approximately 95.81 percent, indicating that a significant amount of variance is still retained after reduction. It is critical that the percentage of explained variance be substantial in order for data reduced to still be a precise representation of the dataset.

You can find my code [here](#).

4 Conclusion

Dimensionality reduction is a significant technique in machine learning with numerous applications. It alleviates the curse of dimensionality by eliminating redundant or noisy features, increasing computing efficiency, and reducing overfitting. Data visualization, feature selection, picture and audio recognition, anomaly detection, and clustering are among its applications. However, it is vital to choose the best technique for the job based on the problem and data characteristics, as well as extensively examine its performance in the application context. There are other difficulties with procedure selection, potential information loss, and identifying the suitably reduced dimensionality. Despite these difficulties, dimensionality reduction is a significant and widely used approach, with continuing research and development aimed at improving its effectiveness and addressing its limits. Researchers and practitioners can acquire useful insights, improve computational efficiency, and overcome the problems associated with high-dimensional datasets by efficiently implementing dimensionality reduction.

5 Declaration of Originality

I, Doluwamu Taiwo Kuye, herewith declare that I have composed the present paper and work by myself and without the use of any other than the cited sources and aids. Sentences or parts of sentences quoted literally are marked as such; other references with regard to the statement and scope are indicated by full details of the publications concerned. The paper and work in the same or similar form have not been submitted to any examination body and have not been published. This paper was not yet, even in part, used in another examination or as a course performance. I agree that my work may be checked by a plagiarism checker.

11/05/2023&Wuppertal - Doluwamu Taiwo Kuye

Bibliography

- [AW10] Abdi, Hervé; Williams, Lynne J: Principal component analysis. Wiley interdisciplinary reviews: computational statistics, 2(4):433–459, 2010. The article presents an in-depth look at PCA, including its mathematical foundations, implementation details, and practical applications.
- [BN06] Bishop, Christopher M; Nasrabadi, Nasser M: Pattern recognition and machine learning, volume 4. Springer, 2006. The book provides a thorough introduction to pattern recognition and machine learning. The book covers a wide range of topics, including probability theory fundamentals, supervised and unsupervised learning, Bayesian approaches, and more advanced topics like neural networks and graphical models.
- [Gé22] Géron, Aurélien: Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow. O'Reilly Media, Inc.", 2022. This book provides a hands-on introduction to machine learning with the Python packages Scikit-Learn and TensorFlow. Data preparation, feature engineering, model selection, deep learning, and reinforcement learning are among the topics covered in the book.
- [Ha09] Hastie, Trevor; Tibshirani, Robert; Friedman, Jerome H; Friedman, Jerome H: The elements of statistical learning: data mining, inference, and prediction, volume 2. Springer, 2009. The book covers statistical learning theory and techniques in depth, including linear regression, tree-based algorithms, neural networks, and support vector machines, it also covers dimensionality reduction components. Model selection, regularization, and unsupervised learning are also covered.
- [LH08] LJPvd, Maaten; Hinton, GE: Visualizing high-dimensional data using t-SNE. J Mach Learn Res, 9(2579-2605):9, 2008. This article describes t-SNE (t-distributed stochastic neighbor embedding), a nonlinear dimensionality reduction approach commonly used for visualizing high-dimensional data in two or three dimensions.
- [LS99] Lee, Daniel D; Seung, H Sebastian: Learning the parts of objects by non-negative matrix factorization. Nature, 401(6755):788–791, 1999. This article describes the non-negative matrix factorization (NMF) machine learning algorithm, which may be used to learn the parts of things unsupervised. The authors show how to use NMF to breakdown photographs of faces into its component elements, such as eyes, noses, and mouths.

-
- [MZ13] Ma, Yanyuan; Zhu, Liping: A review on dimension reduction. *International Statistical Review*, 81(1):134–150, 2013. The paper explores numerous approaches and methods for identifying adequate dimensionality in dimension reduction techniques.
- [Pi23] Pinecone.io: , Dimensionality Reduction. <https://www.pinecone.io/learn/dimensionality-reduction/>, Accessed 2023. Gives a well-defined explanation on dimensionality reduction, also has relevant diagrams.
- [R23] R, ELAVARASAN: , A Quick Overview of Machine Learning Tasks. <https://www.c-sharpcorner.com/article/a-quick-overview-of-machine-learning-tasks/>, Accessed 2023. This article covers different techniques of machine learning.
- [RM17] Raschka, Sebastian; Mirjalili, Vahid: *Python Machine Learning*. Packt Publishing, 2nd edition, 2017. The book discusses a wide range of machine learning topics, including as supervised and unsupervised learning, deep learning, natural language processing, and data visualization. It also contains an introduction to Python for machine learning.
- [Va09] Van Der Maaten, Laurens; Postma, Eric; Van den Herik, Jaap et al.: Dimensionality reduction: a comparative. *J Mach Learn Res*, 10(66-71):13, 2009. The article focuses on dimensionality reduction and provides a comparison of several dimensionality reduction strategies. The writers compare various methods in terms of performance, strengths, and limitations.