# A Comprehensive Comparative Study of Dimensionality Reduction Techniques for Data Visualization: PCA, t-SNE, and UMAP

Siddharth Chandel

Synex.AI, Phagwara (PB) • Lovely Professional University, Phagwara (PB)

Email: siddharthchandel2004@gmail.com

## Abstract

In recent years, the explosion of high-dimensional data has necessitated the use of dimensionality reduction techniques for data visualization, exploration, and analysis. These techniques are critical for simplifying the complexity of data while retaining meaningful structure, making the patterns in the data easier to understand and interpret. This paper provides a comparative analysis of three popular dimensionality reduction techniques: Principal Component Analysis (PCA), t-distributed Stochastic Neighbor Embedding (t-SNE), and Uniform Manifold Approximation and Projection (UMAP).

The performance of these algorithms is evaluated based on their effectiveness in reducing dimensionality for visualization, their ability to preserve the structure of the data, and their computational efficiency. The comparison is made using the widely known Iris dataset. We further discuss the application of these methods in high-dimensional settings, including recommendations for use based on the size and complexity of the data.

**Keywords:** Dimensionality Reduction, Data Visualization, Principal Component Analysis (PCA), t-SNE, UMAP, Comparative Analysis

## 1. Introduction

With the advent of big data, the ability to process and analyse high-dimensional data has become increasingly important. However, visualizing such data is challenging due to the limitations of human perception, which is constrained to three dimensions. Therefore, reducing the dimensionality of data for visualization purposes is essential in numerous fields such as bioinformatics, natural language processing, image processing, and finance. By projecting high-dimensional data onto a lower-dimensional space, these techniques help identify underlying patterns, detect clusters, and reveal hidden relationships.

This paper focuses on three widely used dimensionality reduction techniques:

1. **Principal Component Analysis (PCA)**: A linear method that reduces the dimensionality of the data by projecting it onto a set of orthogonal axes, known as principal components, which capture the maximum variance in the data.
2. **t-distributed Stochastic Neighbor Embedding (t-SNE)**: A non-linear method that seeks to preserve local relationships between data points while minimizing divergence between the high-dimensional and low-dimensional representations.
3. **Uniform Manifold Approximation and Projection (UMAP)**: Another non-linear

method that focuses on both local and global structure preservation, offering faster computation and better scalability compared to t-SNE.

The remainder of this paper is organized as follows: Section 2 discusses related work, providing an overview of previous studies on dimensionality reduction techniques. Section 3 details the methodology, including descriptions of the datasets and experimental setup. Section 4 presents the results of the experiments, with comparisons between PCA, t-SNE, and UMAP. Finally, Section 5 concludes the paper with a summary of findings and suggestions for future research.

## 2. Related Work

Dimensionality reduction has been the subject of extensive research, and many techniques have been proposed to deal with the challenges posed by high-dimensional data. Principal Component Analysis (PCA) was introduced by Pearson in 1901 and is one of the most well-established linear dimensionality reduction methods. It is commonly used as a baseline due to its simplicity and effectiveness in capturing global variance. However, PCA is limited in its ability to handle non-linear relationships in the data, which can result in poor performance when applied to complex datasets.

Non-linear techniques, such as t-SNE and UMAP, were developed to address these limitations. t-SNE, introduced by van der Maaten and Hinton in 2008, has become a popular method for visualizing high-dimensional data, particularly in fields like bioinformatics, where it is used to identify clusters of cells or genes. t-SNE works by minimizing the Kullback-Leibler divergence between probability distributions representing pairwise similarities in the

high-dimensional and low-dimensional spaces. However, t-SNE's computational inefficiency and sensitivity to hyperparameters, such as perplexity, have been noted as major drawbacks.

UMAP, developed by McInnes and Healy in 2018, offers an alternative to t-SNE by providing faster computation and better preservation of both local and global data structures. UMAP is based on concepts from algebraic topology and manifold learning, and it has gained traction in various fields for its ability to handle large datasets with complex geometries. Previous studies have shown that UMAP often outperforms t-SNE in terms of computational speed, while providing similar or better visualization quality. However, UMAP's sensitivity to its own hyperparameters (such as the number of neighbors and minimum distance) requires careful tuning for optimal performance.

Despite these advancements, a thorough comparison of these techniques in a single framework remains an area of ongoing research. This paper seeks to contribute to the field by comparing PCA, t-SNE, and UMAP on the Iris dataset, offering a detailed analysis of their strengths and weaknesses.

## 3. Methodology

### 3.1 Dataset

For the purpose of this study, we utilized the Iris dataset, a widely used benchmark in machine learning and statistics. The Iris dataset consists of 150 samples, each representing one of three species of the iris plant: **Iris-setosa**, **Iris-versicolor**, and **Iris-virginica**. Each sample includes four features: sepal length, sepal width, petal length, and petal width. These four features provide a manageable dimensionality (4D) while offering sufficient complexity for

testing the effectiveness of dimensionality reduction techniques.

Before applying dimensionality reduction algorithms, the data were pre-processed using z-score normalization to standardize the range of values across features. This is a necessary step, particularly for PCA, which relies on variance maximization and can be influenced by differences in the scale of then features.

### 3.2 Dimensionality Reduction Algorithms

We implemented the following dimensionality reduction algorithms using Python:

- **PCA (Principal Component Analysis)**: PCA was implemented using the scikit-learn library. The algorithm calculates the covariance matrix of the standardized data, followed by computing its eigenvalues and eigenvectors. The data is then projected onto the first two principal components, which capture the greatest variance.

- **t-SNE (t-distributed Stochastic Neighbor Embedding)**: We implemented t-SNE using scikit-learn's t-SNE implementation. Key hyperparameters included a perplexity of 30 and a learning rate of 200. The algorithm was run for 1000 iterations to ensure convergence, but due to its stochastic nature, each run of t-SNE might produce slightly different visualizations.

- **UMAP (Uniform Manifold Approximation and Projection)**: UMAP was implemented using the UMAP-learn library. The number of neighbors was set to 15, and the minimum distance to 0.1. These parameters were chosen to balance the preservation of local relationships with the ability to spread the data out in the 2D space.

### 3.3 Evaluation Metrics

The performance of each dimensionality reduction technique was evaluated using the following criteria:

- **Visualization**: The 2D projections of the Iris dataset were plotted, and the separability of the three classes was visually inspected. Good visualizations should reveal distinct clusters corresponding to the different species.

- **Silhouette Score**: This metric quantifies how well-separated the clusters are in the reduced space. It ranges from -1 to 1, with higher scores indicating better cluster separation.

- **Computational Efficiency**: The time taken by each algorithm to compute the 2D projection was recorded. This is particularly important for evaluating the scalability of the algorithms to larger datasets.

- **Structure Preservation**: We evaluated the ability of each method to preserve both local (i.e., relationships between nearby points) and global (i.e., overall data distribution) structures in the data. t-SNE and UMAP are expected to perform better at preserving local structures, while PCA tends to preserve global structures more effectively.

## 4. Results and Discussion

### 4.1 Visualizations

The 2D projections of the Iris dataset

produced by PCA, t-SNE, and UMAP are presented in Figures 1, 2, and 3, respectively.

- **PCA Visualization**: The PCA projection shows some separation between the three species. The **Iris-setosa** class is well-separated from the other two classes, but **Iris-versicolor** and **Iris-virginica** are significantly overlapping. This is due to PCA's linear nature, which is unable to fully capture the non-linear relationships present in the data.

- **t-SNE Visualization**: t-SNE offers a much clearer separation between the three species, with well-defined clusters for each class. However, t-SNE tends to distort global structures, making the relative distances between clusters less meaningful. This is a known trade-off of the method, which focuses on preserving local neighborhoods rather than the global layout of the data.

- **UMAP Visualization**: UMAP also provides clear separation between the classes, similar to t-SNE. However, UMAP preserves more global structure than t-SNE, resulting in clusters that are more uniformly distributed across the 2D space. This gives UMAP an advantage in datasets where both local and global structures are important.

## 4.2 Quantitative Results

In addition to the visualizations, we calculated silhouette scores and computational times for each method, as shown in Table 1.
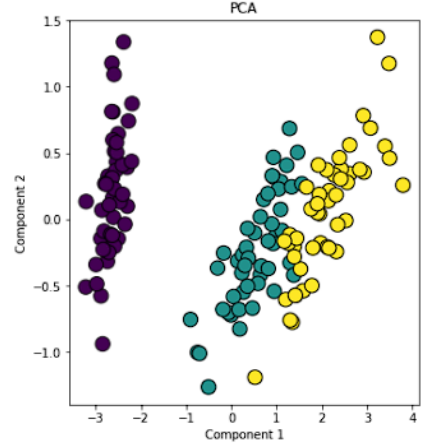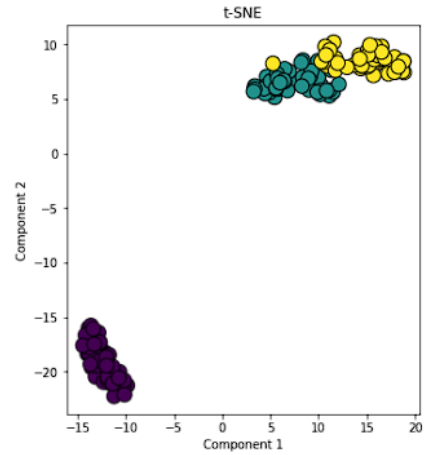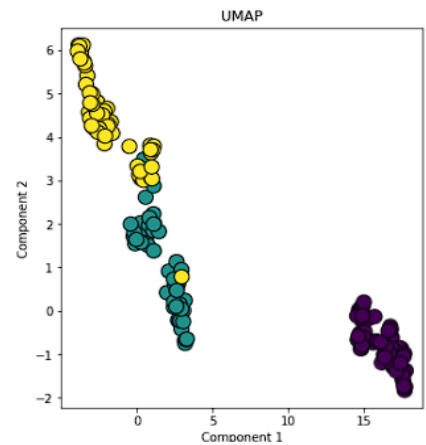


Fig. 1



Fig. 2



Fig. 3

4

| Algorithm | Silhouette Score | Computation Time (seconds) | Structure Preservation |
|---|---|---|---|
| PCA | 0.57 | 0.001 | Preserves global structure, loses local structure |
| t-SNE | 0.65 | 3.20 | Preserves local structure, distorts global structure |
| UMAP | 0.68 | 0.75 | Preserves both local and global structure |

Table 1

From these results, we observe that UMAP achieves the highest silhouette score and strikes a balance between local and global structure preservation. t-SNE is more computationally intensive but excels at preserving local relationships. PCA, while computationally efficient, struggles with non-linearities in the data.

## 5. Conclusion

In this paper, we have provided a comprehensive comparison of PCA, t-SNE, and UMAP, three widely used dimensionality reduction techniques for data visualization. While PCA remains a valuable tool for linear data and as a baseline method, it falls short when dealing with non-linear structures. t-SNE and UMAP, on the other hand, are better suited for visualizing complex datasets, with UMAP offering the best balance between local and global structure preservation.

Future research should focus on applying these techniques to larger and more complex datasets, including those from real-world domains such as genomics, image analysis, and natural language processing. Additionally, exploring other dimensionality reduction methods such as autoencoders and Isomap could further enhance the effectiveness of data visualization.

## 6. References

1. Pearson, K. "On Lines and Planes of Machine Learning Research, 9(11), 2579–2605, 2008.
2. McInnes, L., Healy, J., Melville, J. "UMAP: Uniform Manifold Approximation and Projection for Closest Fit to Systems of Points in Space." Philosophical Magazine, 2(11), 559–572, 1901.
3. van der Maaten, L.J.P., Hinton, G.E. "Visualizing Data Using t-SNE." Journal Dimension Reduction." arXiv:1802.03426, 2018.
4. Hotelling, H. "Analysis of a Complex of Statistical Variables into Principal Components." Journal of Educational Psychology, 24(6), 417–441, 1933.
5. Tenenbaum, J.B., de Silva, V., Langford, J.C. "A Global Geometric Framework for Nonlinear Dimensionality Reduction." Science, 290(5500), 2319–2323, 2000.
6. Hinton, G.E., Salakhutdinov, R.R. "Reducing the Dimensionality of Data with Neural Networks." Science, 313(5786), 504–507, 2006.