# Dimensionality Reduction: A Key Technique in Machine Learning

Sri Poojitha Mandali, Rasika Kole, Vasudha Ambre

October 2024

## 1 Introduction

As the volume and complexity of data continue to grow in the field of data science, it becomes crucial to adopt effective techniques to handle, analyze, and interpret high-dimensional datasets. One of the most powerful methods to achieve this is *Dimensionality Reduction*—a process that reduces the number of input variables while preserving the essential patterns and information contained within the dataset. In this chapter, we will explore how dimensionality reduction techniques can be applied to real-world data, with a particular focus on their role in improving data analysis, visualization, and model efficiency in machine learning.

### 1.1 Research Question

Our central research question is: *How does the inclusion of Spotify track features enhance data analysis, visualization, and model efficiency through the application of dimensionality reduction techniques, such as PCA and t-SNE?*

### 1.2 Why the Question is Interesting and Relevant

In today's data-driven world, datasets with high-dimensional features are becoming the norm across various fields, including music streaming services like Spotify. However, not all features equally contribute to meaningful insights or predictive models. Often, redundant or irrelevant features can cloud the analysis and lead to less efficient machine learning models.

By applying dimensionality reduction techniques, we can address several key issues:

- **Improved Efficiency**: Reducing the number of features in the dataset minimizes memory usage and processing time, making it more feasible to handle large datasets and run machine learning models.

- **Enhanced Model Performance**: By focusing on the most important features, dimensionality reduction helps improve model generalization, reducing overfitting and boosting accuracy.

- **Better Visualization**: High-dimensional data is inherently difficult to visualize. Techniques like PCA and t-SNE allow for meaningful visualization by projecting the data into two or three dimensions, enabling us to uncover patterns, clusters, or trends that are otherwise difficult to observe.

- **Real-World Applications**: In industries like music, finance, and healthcare, vast amounts of data are generated daily. Dimensionality reduction allows data scientists to extract valuable insights, enabling businesses to make data-driven decisions more efficiently.

Ultimately, dimensionality reduction is a vital tool in the machine learning toolkit. It not only improves computational efficiency but also helps maintain the interpretability of complex datasets, leading to more robust and actionable models.

# 2    Theory and Background

Dimensionality reduction is a powerful approach used in data science to tackle high-dimensional datasets efficiently. It serves two major purposes: reducing computational complexity and improving the generalization of machine learning models. Dimensionality reduction can be broadly categorized into two main approaches: **Feature Selection** and **Feature Extraction**. Each of these techniques serves a specific purpose and can be applied depending on the nature of the dataset and the problem being solved.

## 2.1    Feature Selection Techniques

Feature selection focuses on identifying and retaining the most important features from the original dataset without transforming them into a new space. This is particularly useful when some features are highly redundant or irrelevant, which can lead to overfitting or unnecessarily increased model complexity. The goal is to maintain the predictive power of the dataset while reducing its size. Feature selection techniques are generally classified into three types:

### 2.1.1    Filter Methods

Filter methods are relatively simple, computationally inexpensive techniques that rank features based on statistical measures such as correlation, variance, or chi-squared tests. Features are selected based on these scores, and only the top-ranking ones are retained. While these methods are fast and easy to implement, they do not account for interactions between features. Common filter methods include correlation coefficients and mutual information.

### 2.1.2 Wrapper Methods

Wrapper methods involve evaluating subsets of features by training a machine learning model on each subset. The model's performance determines which features are important. A popular wrapper technique is **Recursive Feature Elimination (RFE)**, which recursively removes the least important features until a desired number of features is reached. While wrapper methods typically result in higher model performance compared to filter methods, they are more computationally expensive since they require multiple iterations of model training.

### 2.1.3 Embedded Methods

Embedded methods perform feature selection during the model training process itself. Unlike filter and wrapper methods, which are independent of the model, embedded methods consider feature importance while building the model. One well-known example is **Lasso Regression**, where a regularization penalty is applied to the coefficients of less important features, effectively shrinking some of them to zero. This process both simplifies the model and selects relevant features.

## 2.2 Feature Extraction Techniques

Feature extraction goes beyond selecting existing features; it involves transforming the data into a new feature space where the most important characteristics of the data are captured. This transformation is particularly useful when the original features are highly correlated or not well-suited for direct analysis. Feature extraction techniques include both linear and non-linear models:

### 2.2.1 Principal Component Analysis (PCA)

PCA is a widely used linear dimensionality reduction technique. It transforms the dataset into a new coordinate system where the axes (called principal components) correspond to directions of maximum variance in the data. The first principal component captures the highest variance, followed by the second, and so on. PCA reduces dimensionality by projecting the data onto the first few principal components, effectively reducing the number of features while retaining the majority of the variance.

**Steps Involved in PCA:**

1. Standardize the dataset so that each feature has a mean of 0 and variance of 1.

2. Compute the covariance matrix to understand how the features vary in relation to each other.

3. Calculate the eigenvalues and eigenvectors of the covariance matrix. The eigenvectors represent the directions of maximum variance, while the eigenvalues represent the magnitude of that variance.

3

4. Sort the eigenvectors by their eigenvalues in descending order.

5. Select the top eigenvectors (principal components) and project the data onto this lower-dimensional space.

PCA is effective for reducing dimensionality while retaining much of the data's inherent variability, making it ideal for tasks like noise reduction and visualization.

### 2.2.2   t-Distributed Stochastic Neighbor Embedding (t-SNE)

t-SNE is a non-linear technique designed for dimensionality reduction, particularly useful for visualizing high-dimensional data. Unlike PCA, which preserves global variance, t-SNE focuses on maintaining local structure. It does this by minimizing the divergence between two probability distributions: one that measures pairwise similarities in the high-dimensional space and another that measures similarities in the low-dimensional space. This makes t-SNE an excellent tool for clustering and visualizing data.

**Steps Involved in t-SNE:**

1. Compute pairwise similarities between data points in the high-dimensional space.

2. Convert these similarities into conditional probabilities representing the likelihood that one point would be a neighbor of another.

3. Compute similar probabilities in the lower-dimensional space, aiming to maintain local similarities.

4. Minimize the Kullback-Leibler divergence between the two probability distributions using gradient descent.

5. The resulting lower-dimensional representation captures the local structure of the data, allowing for effective visualization.

t-SNE is particularly effective for uncovering hidden patterns, such as clusters, in high-dimensional data, making it popular in fields like genomics and image recognition.

# 3   Problem Statement

## 3.1   Detailed Problem Statement:

The Spotify Tracks dataset contains a large variety of numerical and categorical features that describe various aspects of songs, such as danceability, energy, tempo, and loudness. With so many features, the dataset becomes difficult to analyze, visualize, and interpret due to its high dimensionality. High-dimensional data can lead to overfitting, increased computational complexity, and challenges in identifying meaningful patterns.

To address these issues, we employ a combination of feature selection techniques, such as filter methods, wrapper methods, and embedded methods, and feature extraction techniques, including Principal Component Analysis (PCA)and t-Distributed Stochastic Neighbor Embedding (t-SNE). The goal of these techniques is to reduce the dimensionality of the dataset while retaining crucial information and patterns. This process will help simplify the dataset, making it easier to handle and more suitable for machine learning models.

## 3.2 Input-Output Format:

**Input:** The input is a high-dimensional dataset of Spotify tracks, represented as $X$. The dataset contains a variety of features, such as:

- Danceability

- Tempo

- Loudness

- Energy

- Duration

- Acousticness

- Liveness

**Output:** The output is a transformed dataset, represented as $X'$, which has been reduced in dimensionality through one or more of the following techniques: feature selection (filter, wrapper, or embedded methods) or feature extraction (PCA, t-SNE, Autoencoders, or Isomap). The reduced dataset retains the most significant patterns and information from the original data, while having fewer dimensions for easier analysis and visualization.

## 3.3 Sample Input-Output:

Consider a dataset where each track has 15 different features. Using one of the feature extraction methods, such as PCA, the number of features is reduced to 3 principal components, while preserving 90% of the variance in the data. This reduced dataset contains most of the important information from the original dataset but is now much easier to work with.

**Sample Input:**

- Danceability: 0.78

- Energy: 0.85

- Loudness: -6.0 dB

- Tempo: 125 BPM

- Duration: 210000 ms

- Valence: 0.90

- Acousticness: 0.15

- Liveness: 0.30

- Speechiness: 0.05

- Instrumentalness: 0.10

**Sample Output (PCA Reduction to 2 Components):**

- Principal Component 1: 1.25

- Principal Component 2: -0.75

**Sample Output (t-SNE Reduction to 2 Dimensions):**

- t-SNE Dimension 1: -2.35

- t-SNE Dimension 2: 0.80

By applying these dimensionality reduction techniques, the dataset is simplified while retaining its most important characteristics, allowing for easier handling, analysis, and visualization. The resulting lower-dimensional dataset can be used in further machine learning models, improving both efficiency and interpretability.

# 4 Problem Analysis

## 4.1 Constraints

**Interpretability:** In dimensionality reduction, the balance between reducing complexity and maintaining interpretability is a critical constraint. In this project, we applied both **feature selection** and **feature extraction** techniques to reduce the dimensionality of the dataset. Feature selection methods, such as filter methods, wrapper methods, and embedded methods, are effective in identifying and removing irrelevant or redundant features. However, the interpretability of these selected features can sometimes be limited due to interactions between the features themselves.

Feature extraction techniques, such as PCA and t-SNE, further reduce dimensionality by transforming the data into a new feature space, often at the cost of interpretability. PCA, for example, combines features into principal components, making it difficult to interpret these new components as they are often a combination of multiple original features. While t-SNE excels in visualizing the local structure of the data, it is less interpretable since it does not produce a linear transformation like PCA.
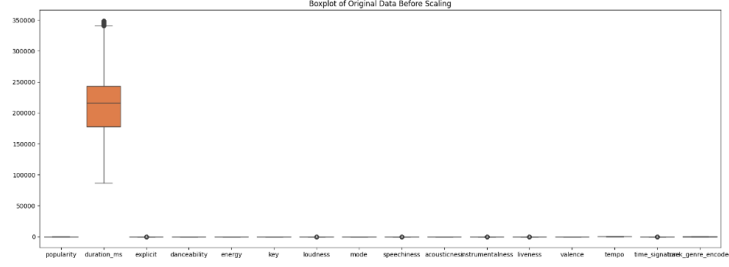
Figure 1: Distribution of different audio features in the Spotify dataset before scaling, highlighting the disparities in their ranges.

**Computational Complexity:** When dealing with a high-dimensional dataset like the Spotify Tracks dataset, computational complexity is a major consideration. Feature selection techniques, especially filter methods, are computationally efficient as they rank features based on statistical criteria without needing to train models. However, wrapper methods such as Recursive Feature Elimination (RFE) are computationally more expensive since they involve model training and evaluation for each subset of features.

Similarly, t-SNE is computationally intensive because it calculates the pairwise similarities between data points. To handle this, PCA is used as a preprocessing step, reducing the dimensionality of the data before applying t-SNE, which limits the computational burden.

## 4.2   Approach to Solve the Problem

To simplify the high-dimensional Spotify dataset and make it more suitable for analysis and visualization, a combination of feature selection and feature extraction techniques was applied. The process can be broken down into two main steps:

1. Feature Selection: We began by applying feature selection techniques to identify the most relevant features in the dataset. Filter methods, such as correlation ranking, were used to evaluate the relationship between features and the target variable. Additionally, wrapper methods like Recursive Feature Elimination (RFE) were applied to refine the feature set by iteratively training models and removing the least important features. Embedded methods, like Lasso Regression, were also employed to select features during model training by penalizing less important features.

2. Feature Extraction: Once the most relevant features were selected, we applied feature extraction techniques to further reduce the dimensionality. PCA was used to capture the maximum variance in the dataset while reducing the feature space to a manageable size. Finally, t-SNE was applied to visualize the relationships and clusters in the data, particularly

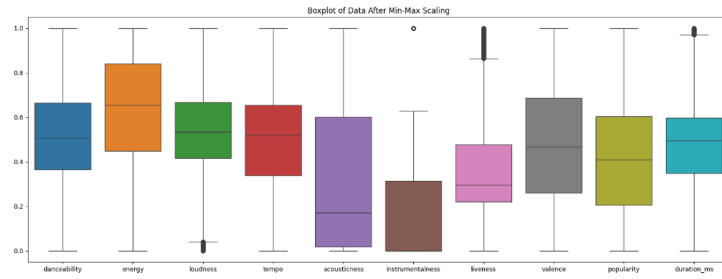focusing on local similarities between tracks.



Figure 2: Boxplot of features after scaling, ensuring that all features contribute equally to the PCA process by bringing them to a common range between 0 and 1.

This logical approach allows us to reduce the dataset's complexity through feature selection while maintaining the core information through feature extraction. By reducing the dimensions in this step-by-step process, we enable more effective data visualization and machine learning model training.

## 4.3 Identification of Key Data Science and Algorithmic Principles

**Feature Selection Techniques:** The key principle behind feature selection is identifying the most relevant features for a given task. Filter methods use statistical measures like correlation or chi-squared tests to rank features, while wrapper methods evaluate feature subsets based on model performance. Embedded methods, such as Lasso Regression, combine feature selection with model training, applying penalties to reduce the importance of less relevant features.

**Variance Retention in PCA:** PCA is based on the principle of variance retention. By identifying the directions in which the data exhibits the most variance, PCA reduces the number of dimensions while preserving the most critical information. This ensures that the reduced dataset maintains its key patterns and relationships, allowing for more efficient data analysis.

**Local Relationship Preservation in t-SNE:** Unlike PCA, which retains global variance, t-SNE focuses on preserving local relationships between data points. It maps pairwise similarities into a lower-dimensional space, making it useful for identifying clusters and patterns, such as grouping songs with similar audio characteristics in the Spotify dataset.

The correlation matrix, for example, shows the relationships between different features in the Spotify dataset, highlighting the importance of reducing redundancy through feature selection and extraction techniques.
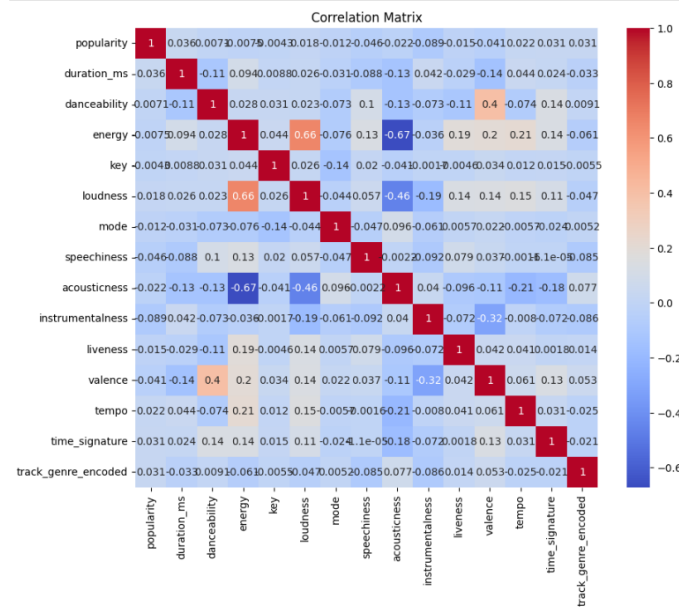
Figure 3: Correlation matrix of numerical features in the dataset. This helps identify relationships between features and highlights redundancy, which is addressed through dimensionality reduction.

# 5 Solution

This section describes the step-by-step approach used to reduce the dimensionality of the Spotify Tracks dataset, from identifying and selecting relevant features to visualizing the reduced feature space. The workflow is designed to systematically select the most significant features while maintaining interpretability and computational efficiency.

## 5.1 Step-by-Step Solution

The solution can be broken down into the following steps:

**Step 1: Feature Selection using Wrapper and Embedded Methods**
We first applied feature selection techniques to reduce the number of irrelevant or redundant features. These techniques help in retaining only the most important features that significantly impact the target variable.

- **Wrapper Methods (Recursive Feature Elimination - RFE)**: RFE recursively removes the least important features and ranks them based on their importance to the model. This method ensures that only the most relevant features remain after each iteration. We used this method to evaluate the impact of each feature on the model's predictive power.

- **Embedded Methods (Lasso Regression):** Lasso regression adds a regularization term that shrinks the coefficients of less important features to zero. This not only selects important features during model training but also prevents overfitting by penalizing large coefficients. By using Lasso, we further reduced the number of features while maintaining a robust model.

**Step 2: Further Feature Evaluation with OLS Regression and Permutation Importance** Once the most relevant features were selected using RFE and Lasso, we applied **OLS Regression (Ordinary Least Squares)** to analyze the linear relationships between the selected features and the target variable. This regression technique helped us better understand how each feature contributes to the prediction.

In addition, **Permutation Importance** was used to assess the impact of individual features by randomly shuffling each feature and observing the decrease in model performance. This method helps in validating the significance of the remaining features and ensures that only the most relevant ones are retained for further analysis.
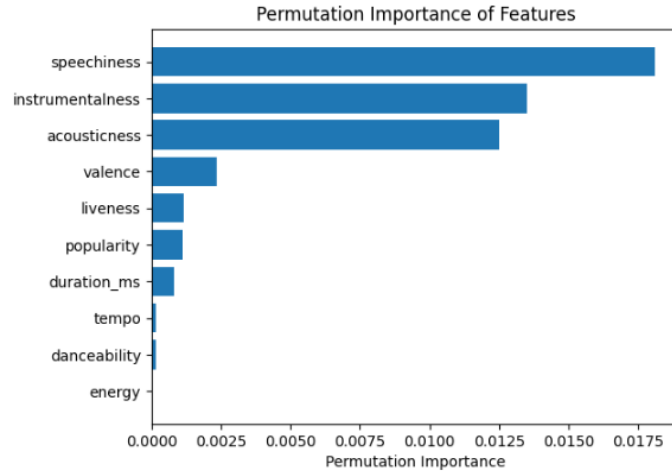


Figure 4: Feature ranking based on Permutation Importance, showing the impact of each feature on the model's performance.

**Step 3: Dimensionality Reduction using PCA (Principal Component Analysis)** After identifying the most significant features, we applied **PCA** to reduce the feature space further while retaining as much variance as possible. PCA transforms the dataset into a new coordinate system, where each axis represents a principal component that captures the highest variance in the data.

**Step 4: Visualization with t-SNE (t-Distributed Stochastic Neighbor Embedding)**   Finally, we applied \*\*t-SNE\*\* to visualize the dataset in a two-dimensional space.  t-SNE is a non-linear dimensionality reduction technique that focuses on preserving the local relationships between data points, making it ideal for visualizing clusters and patterns within the dataset.

## 5.2   Pseudocode:

```
# Step 1: Standardize the Data
from sklearn.preprocessing import MinMaxScaler
scaled_data = MinMaxScaler().fit_transform(original_data)

# Step 2: Apply Recursive Feature Elimination (RFE) for Feature Selection
from sklearn.feature_selection import RFE
from sklearn.linear_model import LinearRegression
model = LinearRegression()
rfe = RFE(estimator=model, n_features_to_select=10)
rfe_data = rfe.fit_transform(scaled_data, target_variable)

# Step 3: Apply Lasso for Embedded Feature Selection
from sklearn.linear_model import Lasso
lasso = Lasso(alpha=0.01)
lasso.fit(rfe_data, target_variable)
important_features = lasso.coef_ != 0

# Step 4: Apply PCA for Dimensionality Reduction
from sklearn.decomposition import PCA
pca = PCA(n_components=3)
pca_data = pca.fit_transform(rfe_data[:, important_features])

# Step 5: Apply t-SNE for Visualization
from sklearn.manifold import TSNE
tsne = TSNE(n_components=2)
tsne_results = tsne.fit_transform(pca_data)

# Step 6: Plot the t-SNE Results
import matplotlib.pyplot as plt
plt.scatter(tsne_results[:, 0], tsne_results[:, 1], c=genres, cmap='viridis')
plt.xlabel('Component 1')
plt.ylabel('Component 2')
plt.title('t-SNE Visualization of Spotify Data')
plt.show()
```

## 5.3 Proof of Correctness

The correctness of this approach is supported by the logical sequence of steps. First, feature selection techniques like RFE and Lasso were used to ensure that only the most relevant features were retained. This eliminated noise and redundant information from the dataset. OLS Regression and Permutation Importance provided further validation by quantifying the impact of individual features.

Next, PCA was applied to retain the maximum variance in the dataset while reducing the dimensionality, ensuring that no critical information was lost. Finally, t-SNE allowed for an intuitive visualization of the data in two dimensions, highlighting clusters and relationships between tracks that were not apparent in the original high-dimensional space.

The combination of these methods not only reduced the dataset's complexity but also improved its interpretability and provided a robust solution for dimensionality reduction.

# 6 Results and Data Analysis

## 6.1 Well-Presented Results with Visualizations:

**PCA Results:** A cumulative variance plot was generated to show the proportion of variance captured by the principal components. The plot indicated that the first three components captured approximately 90% of the variance.

**t-SNE Results:** A 2D scatter plot was created from the t-SNE results, where each point represented a track and was colored by genre.

## 6.2 Insightful Discussion of Results and Their Implications:

**PCA:** By reducing the features from 15 to 3, further processing was made possible to support easy visualization without losing critical information.

**t-SNE:** The clustering in the 2D scatter plot showed that similar tracks were grouped together, which is useful for making a music recommendation system.

# 7 Connection to Theory

**PCA and Variance Maximization:** By focusing on the directions of maximum variance, PCA captures important structure in the data.

**t-SNE and Clustering:** t-SNE conserves local similarities, allowing for intuitive clustering of data, useful for identifying genres or for playlist generation.
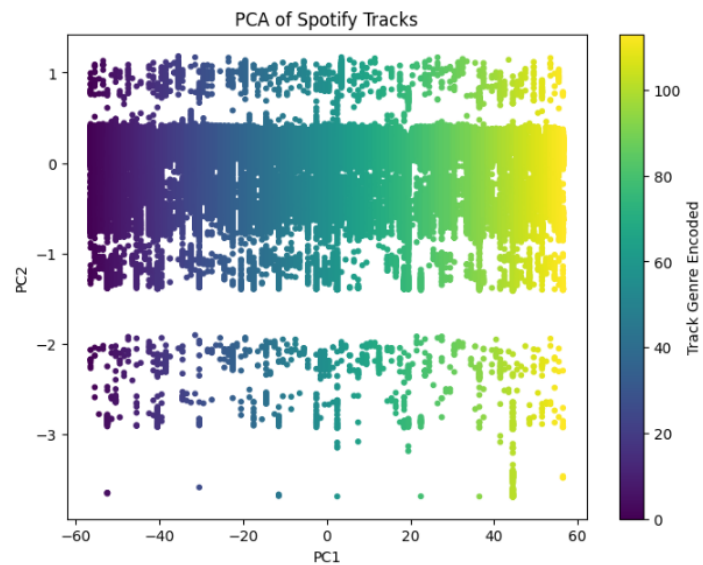
Figure 5: Cumulative variance plot helped determine the optimal number of components to use.
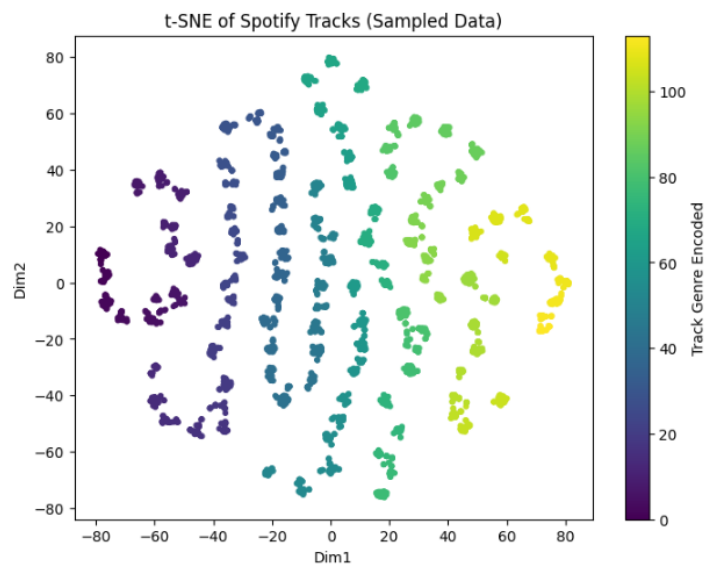


Figure 6: 2D scatter plot of the tracks clustered by similar audio features, colored by genre.

# 8    Conclusion

In this project, we successfully reduced the dimensionality of the Spotify dataset by applying feature selection techniques such as RFE and Lasso, followed by PCA and t-SNE for feature extraction and visualization. These methods allowed us to retain the most important features and capture over 90% of the data's variance while simplifying the dataset for easier interpretation. The t-SNE visualization revealed clear clusters in the data, highlighting meaningful patterns. Overall, this approach not only improved computational efficiency but also enhanced the dataset's interpretability. This dimensionality reduction framework provides a solid foundation for future machine learning models or further exploratory analysis.

# 9    References

- Jolliffe, I. T. (2002). *Principal Component Analysis*, 2nd Edition. Springer Series in Statistics.

- van der Maaten, L., & Hinton, G. (2008). *Visualizing Data using t-SNE*. Journal of Machine Learning Research, 9(Nov), 2579–2605.

- Pearson, K. (1901). *On Lines and Planes of Closest Fit to Systems of Points in Space.* Philosophical Magazine, 2(11), 559-572.

- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction.* Springer.

- van der Maaten, L. (2014). *Accelerating t-SNE using Tree-Based Algorithms.* Journal of Machine Learning Research, 15(93), 3221-3245.

- Hinton, G. E., & Salakhutdinov, R. R. (2006). *Reducing the Dimensionality of Data with Neural Networks.* Science, 313(5786), 504–507.