**Title:** Dimensionality Reduction for Visualization: Apply techniques like PCA (Principal Component Analysis) to reduce the dimensionality of data for visualization purposes.

**Abstract:**

In today's data-driven fashion industry, large and complex datasets are collected from various sources, including customer behavior, style preferences, and seasonal trends . Analyzing and extracting meaningful insights from such high-dimensional data presents a significant challenge. Dimension reduction techniques like Principal Component Analysis (PCA) play a critical role in simplifying these datasets by reducing the number of features while retaining the most informative patterns. This paper explores the application of PCA in the context of fashion data analysis, illustrating how this method can uncover underlying trends, improve customer segmentation, and streamline inventory management.

The methodology involves applying PCA to a representative fashion dataset, such as Fashion-MNIST or an e-commerce purchase history dataset. By transforming the data into a lower-dimensional space, PCA reveals core patterns that would otherwise remain hidden in high-dimensional spaces. For instance, PCA can identify prominent features like seasonal color trends, popular clothing cuts, or texture patterns that drive consumer preferences. Furthermore, in customer segmentation, PCA enables brands to categorize customers into clusters based on a few principal components, facilitating more personalized marketing and targeted recommendations.

# INTRODUCTION

The fashion industry is increasingly data-driven, relying on vast amounts of data to inform decisions on style trends, customer preferences, and inventory strategies. With the rise of e-commerce, social media, and global supply chains, data on fashion preferences and purchasing behaviors is collected continuously from various sources such as online browsing, purchasing history, product reviews, and even real-time trends on social media platforms. These rich datasets contain valuable insights but also pose significant challenges: they are often high-dimensional, containing hundreds or even thousands of variables that capture aspects of individual products (e.g., color, fabric, cut, and brand) as well as customer demographics, purchase behaviors, and seasonality. High-dimensional data tends to be complex, with redundant or irrelevant features that can increase computational costs, obscure key patterns, and complicate data analysis.

To tackle these issues, dimensionality reduction techniques like Principal Component Analysis (PCA) are essential. PCA reduces the complexity of these datasets by transforming the original high-dimensional data into a lower-dimensional space that retains the most informative characteristics. By finding new, uncorrelated variables, or "principal components," PCA maximizes the variance explained by each component, allowing analysts to capture essential patterns without the noise of redundant data. This transformation not only simplifies data interpretation but also makes it

feasible to perform advanced analyses and visualizations. For example, PCA can reveal dominant fashion trends over time, identify clusters of customers with similar preferences, and enable companies to streamline their inventory by focusing on high-demand items. By uncovering hidden structures within complex fashion data, PCA offers insights that are crucial for trend forecasting, personalized recommendations, and data-driven inventory management.

In this paper, we explore the application of PCA in fashion data analysis, demonstrating how this method helps transform high-dimensional datasets into actionable insights. Through practical examples and visualizations, we show how PCA reduces the complexity of fashion datasets, captures essential patterns, and supports various business applications such as trend analysis, customer segmentation, and inventory optimization. The findings illustrate PCA's value in enabling fashion companies to leverage data for more efficient decision-making, competitive strategy, and enhanced consumer experiences.

# Principal Component Analysis (PCA)

Principal Component Analysis (PCA) is a fundamental technique for dimensionality reduction, transforming high-dimensional data into a set of uncorrelated variables called principal components (PCs) that capture the most variance. PCA is sensitive to data scale, so the dataset XXX is commonly standardized before applying the transformation. Standardization of each feature jjj for data point xij is calculated as:

$$z_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

Where uj is the mean of feature j and $\sigma_j$ its standard deviation.

Once the data is standardized, PCA proceeds by calculating the covariance matrix $\Sigma$ of the dataset, summarizing the linear relationships between

features. For an n × p matrix , the covariance matrix $\Sigma$ is given by:

$$\Sigma = \frac{1}{n-1} X^T X$$

where X^T is the transpose of X . Next, to identify the principal components, we perform eigenvalue decomposition on $\Sigma$ , resulting in eigenvalues $\lambda i$ and eigenvectors vi that satisfy:

$$\Sigma v_i = \lambda_i v_i$$

Each eigenvalue $\lambda i$ corresponds to the amount of variance captured by its eigenvector vi , which defines the direction of each principal component. The principal components are projections of the original data onto these eigenvectors, where the i-th principal component PC i is calculated as:

$$\mathrm{PC}_i = X \cdot v_i$$

Each principal component represents a new axis in the transformed feature space, ordered by the variance it captures.

The total variance of the data is the sum of all eigenvalues. The explained variance ratio for each principal component, indicating the proportion of total variance it captures, is given by:

$$\text{Explained Variance Ratio} = \frac{\lambda_i}{\sum_{j=1}^{p} \lambda_j}$$

To determine the cumulative variance captured by the first kkk components, we calculate:

$$\text{Cumulative Explained Variance} = \sum_{i=1}^{k} \frac{\lambda_i}{\sum_{j=1}^{p} \lambda_j}$$

For dimensionality reduction, we select the top k eigenvectors to form the transformation matrix Vk

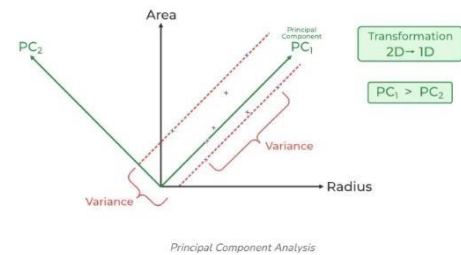, then project the data X onto this reduced k - dimensional space as:

$$Z = X \cdot V_k$$

where Z is the n × k matrix representing the dataset in the lower-dimensional space. These steps—standardization, covariance calculation, eigenvalue decomposition, explained variance analysis, and projection — form the mathematical foundation of PCA, simplifying high-dimensional data and retaining the most informative patterns.

# Applications of PCA in Fashion Data Analysis

PCA has multiple applications in fashion data analysis, enabling businesses to gain valuable insights from complex datasets. In fashion trend analysis, PCA can uncover dominant patterns in colors, cuts, and textures; for example, a primary component might reveal popular colors across seasonal collections, such as brighter hues in summer. For customer segmentation, PCA simplifies customer attributes like age, gender, and purchase history, identifying clusters that reflect core preferences, which allows for targeted marketing. In inventory management, PCA aids in focusing on high-demand items while reducing stock of less popular products, thus optimizing inventory. Finally, PCA supports product recommendations by reducing each item's feature space, allowing for more accurate suggestions based on similar attributes—this is especially useful for visual search tools that match items based on shared characteristics. These applications help fashion companies use data-driven approaches for efficient decision-making across various operational areas.

Principal Component Analysis (PCA) is used to reduce the dimensionality of a data set by finding a new set of variables, smaller than the original set of variables, retaining most of the sample's information, and useful for the regression and classification of data.
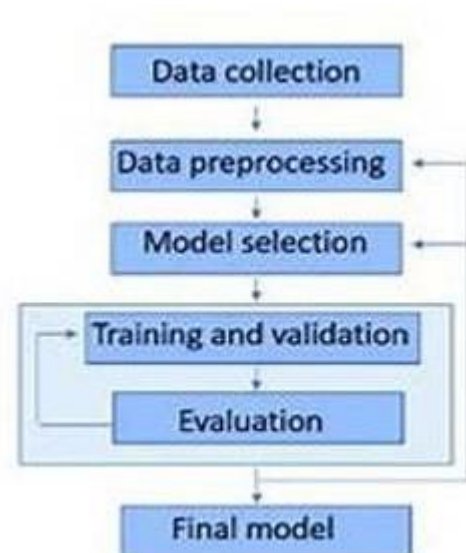


*Principal Component Analysis*

# Methodology: Applying PCA on a Fashion Dataset

**Dataset Description**:

- Choose a dataset (e.g., Fashion-MNIST or customer purchase data from an online retailer).

**Data Preprocessing**:

- Discuss steps like normalization and reshaping images into vectors if image data is used.

- PCA requires centering the data by subtracting the mean from each feature.



**Implementing PCA**:

- **Step 1**: Standardize the data matrix X.

- **Step 2**: Calculate the covariance matrix Σ.

- **Step 3**: Perform eigenvalue decomposition on Σ to find eigenvalues and eigenvectors.

- **Step 4**: Sort eigenvectors by eigenvalues in descending order and select the top k components.

- **Step 5**: Project data onto the new k-dimensional space:

$$Z = X \cdot V_k$$

  where Vk contains the selected top k eigenvectors.

**Tools and Libraries**:

- Utilize scikit-learn, pandas, and matplotlib for PCA implementation.

# Discussion and Results

**PCA (Principal Component Analysis)**:

PCA is used for dimensionality reduction and feature extraction. It transforms the dataset into a new set of orthogonal components that explain the maximum variance in the data.

The code applies PCA with different numbers of components (2, 4, 10, and 120) to analyze the Fashion-MNIST dataset.

**PCA Components and Variance**:

For n_components=2, n_components=4, and n_components=10, PCA components (principal axes) are visualized using heatmaps. Each heatmap corresponds to a principal component reshaped into a 28x28 grid (which represents an image).

The explained variance ratio of each component is also displayed as titles for each heatmap, showing how much variance each component explains.

**Reconstruction of Images**:

The code includes functionality to reconstruct images from the reduced PCA representation

(120 components). The reconstructed images are compared to the original images.

The reconstruction error (difference between the original and reconstructed images) is also visualized.

**Quartile Analysis**:

The function quartile record computes the reconstruction error for each sample and then retrieves the image corresponding to a specific quantile (e.g., 2nd, 10th, 50th, 90th percentile of errors).

The plot quartiles function visualizes the images corresponding to different error percentiles.

**Metric Calculation (Mean Squared Error)**:

Mean squared error (MSE) is used to measure the reconstruction error between the original and reconstructed images. This metric is used for error analysis in the quartile record and record similarity functions.

Histograms are plotted for the reconstruction errors across different PCA components.

**Heatmaps of PCA Components**:
These heatmaps show the first few principal components (each reshaped into 28x28 pixels) and their corresponding explained variance.

**Reconstructed Images**:
Images are reconstructed using a reduced number of PCA components (e.g., 120). The reconstruction is compared to the original images, and the difference (error) is visualized.

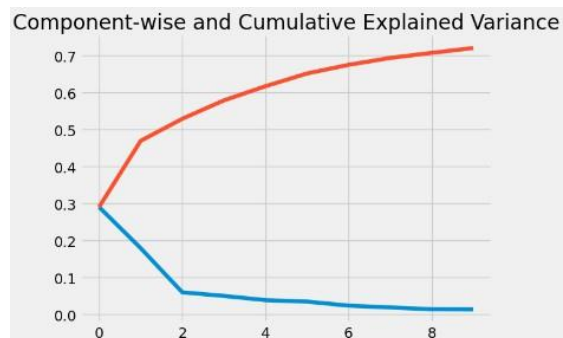**Quartile-Based Reconstruction**:
The error of reconstruction is analyzed at different percentiles, helping to identify which images are harder or easier to reconstruct based on PCA components.

**Histograms of Reconstruction Errors**:
For each PCA component, a histogram of reconstruction errors (MSE) is plotted, showing how well each component can capture the variation in the data
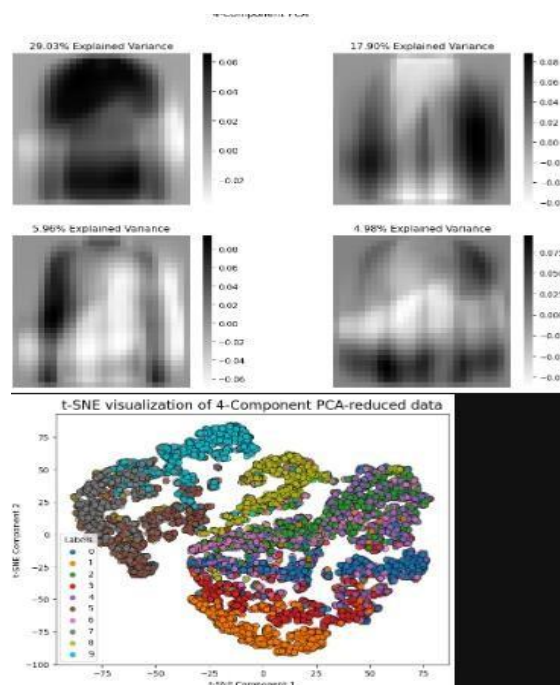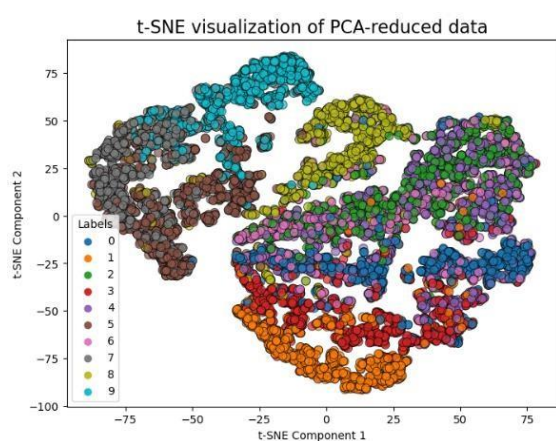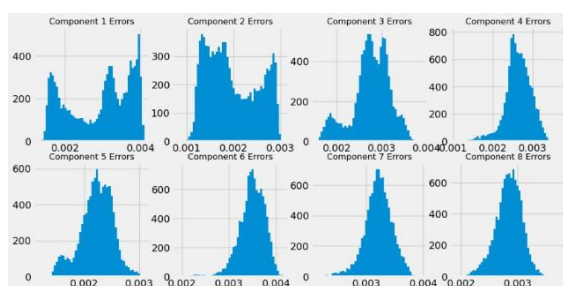
Result :

Component – wise and cumulative Explained Variance

Component-wise and Cumulative Explained Variance



4-Component PCA

## Percentile of fashion clothes





t-SNE visualization of 4-Component PCA-reduced data

Histogram based on result :



^

*Variation in variance : |*

# Advanced PCA Techniques and Alternatives

In addition to standard PCA, several advanced techniques offer enhanced capabilities for handling more complex data structures. **Kernel PCA** is one such extension, which addresses the limitations of linear PCA by mapping data into a higher-dimensional space using kernel functions. This allows Kernel PCA to capture non-linear relationships between features that traditional PCA might miss, making it particularly useful for fashion datasets where complex patterns, such as intricate textures or non-linear customer preferences, exist. On the other hand, **Sparse PCA** introduces a sparsity constraint during the decomposition process, aiming to focus only on the most significant features that have a substantial impact on the principal components. This approach enhances interpretability by reducing the number of features involved in each component, which is valuable when dealing with large, high-dimensional datasets like fashion attributes, as it highlights the most relevant factors driving trends or consumer choices.



t-SNE visualization of PCA-reduced data

When comparing PCA to other dimensionality reduction techniques, **t-SNE** (t-Distributed Stochastic Neighbor Embedding) and **Autoencoders** emerge as alternatives that excel in non-linear data visualization and clustering. While PCA is effective for linear data, **t-SNE** excels in mapping high-dimensional data to a lower-dimensional space, preserving local structure and revealing clusters or patterns that are not linearly separable. This makes it useful in fashion for visualizing customer segmentation or item similarities in an intuitive 2D or 3D space. **Autoencoders**, a type of neural network, also provide a non-linear approach to dimensionality reduction. By learning an efficient encoding of data through a neural network, autoencoders can capture complex, non-linear patterns in fashion datasets, making them ideal for tasks like image recognition, recommendation systems, and anomaly detection in fashion items. Though PCA remains a foundational tool, these advanced methods offer greater flexibility for exploring intricate relationships and providing deeper insights into the fashion industry's dynamic data.

PCA is effective for linear data, **t-SNE** excels in mapping high-dimensional data to a lower-dimensional space, preserving local structure and revealing clusters or patterns that are not linearly separable

# Advantages and Challenges of Using PCA in Fashion Analysis

The use of PCA in fashion analysis offers several advantages. It simplifies data analysis by reducing the dimensionality of complex datasets, making it easier to identify key patterns and trends. By focusing on the principal components that capture the most variance, PCA helps reduce noise and emphasizes the most relevant features. This not only speeds up data processing and analysis but also aids in visualizing high-dimensional data, which can be particularly valuable for uncovering underlying trends in fashion, such as seasonal styles or emerging colour preferences.

However, PCA also presents some challenges. One significant drawback is that it assumes linear relationships between features, potentially discarding subtle, non-linear patterns that might exist in fashion data, such as complex interactions between style elements or customer behavior. Additionally, the transformation to principal components can reduce interpretability, as the original features are combined into new axes that may be difficult to relate directly to specific fashion attributes. This loss of transparency can make it harder for stakeholders to fully understand the reasons behind certain trends or customer preferences, limiting the ability to act on these insights without further analysis or interpretation. Despite these challenges, PCA remains a powerful tool when used appropriately, particularly for uncovering broad patterns in fashion data.

# Conclusion

In summary, Principal Component Analysis (PCA) stands out as an effective tool for reducing the dimensionality of fashion datasets, making complex data more accessible and insightful. By transforming high-dimensional data into a lower-dimensional space while retaining the most significant patterns, PCA simplifies the analysis process and uncovers key trends that may otherwise remain hidden. This approach enables fashion brands and retailers to analyze trends, customer preferences, and product attributes more efficiently, providing a clearer picture of the market dynamics.

Looking towards the future, there are numerous opportunities to integrate PCA with advanced techniques such as deep learning models or autoencoders, especially for fashion image data. Combining PCA with deep learning could enhance the ability to uncover more complex, non-linear relationships in fashion data, while autoencoders could provide an alternative means for dimensionality reduction, particularly in dealing with high-dimensional visual data. These integrations could push the boundaries of what is possible in fashion data analysis, offering deeper and more nuanced insights.

The implications for the fashion industry are substantial. PCA-derived insights can significantly

improve trend forecasting by identifying patterns in customer behavior, seasonal preferences, and product features. Furthermore, PCA helps in segmenting customers based on key preferences, which can guide targeted marketing strategies and personalized product recommendations. By streamlining data analysis, PCA enables companies to make more informed decisions about inventory management, product launches, and customer engagement, ultimately enhancing the overall customer experience and driving business success in the competitive fashion industry.

# REFERENCE

- Ivosev, Gordana, Lyle Burton, and Ron Bonner. "Dimensionality reduction and visualization in principal component analysis." *Analytical chemistry* 80.13 (2008): 4933-4944.
- Engel, Daniel, Lars Hüttenberger, and Bernd Hamann. "A survey of dimension reduction methods for high-dimensional data analysis and visualization." *Visualization of Large and Unstructured Data Sets: Applications in Geospatial Planning, Modeling and Engineering-Proceedings of IRTG 1131 Workshop 2011*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik, 2012.
- Tsai, Flora S. "Dimensionality reduction techniques for blog visualization." *Expert Systems with Applications* 38.3 (2011): 2766-2773.

- Yin, Hujun. "Nonlinear dimensionality reduction and data visualization: a review." *International Journal of Automation and Computing* 4 (2007): 294-303.


- Hasan, Basna Mohammed Salih, and Adnan Mohsin Abdulazeez. "A review of principal component analysis algorithm for dimensionality reduction." *Journal of Soft Computing and Data Mining* 2.1 (2021): 20-30.
- Holbrey, Richard. "Dimension reduction algorithms for data mining and visualization." *Edinburgh: University of Leeds* (2006).
- Nanga, Salifu, et al. "Review of dimension reduction methods." *Journal of Data Analysis and Information Processing* 9.3 (2021): 189-231.
- Sacha, Dominik, et al. "Visual interaction with dimensionality reduction: A structured literature analysis." *IEEE transactions on visualization and computer graphics* 23.1 (2016): 241-250.
- An, Jiyuan, et al. "A dimensionality reduction algorithm and its application for interactive visualization." *Journal of Visual Languages & Computing* 18.1 (2007): 48-70.
- Ayesha, Shaeela, Muhammad Kashif Hanif, and Ramzan Talib. "Overview and comparative study of dimensionality reduction techniques for high dimensional data." *Information Fusion* 59 (2020): 44-58.
- Geng, Xin, De-Chuan Zhan, and Zhi-Hua Zhou. "Supervised

nonlinear dimensionality reduction for visualization and classification." *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 35.6 (2005): 1098-1107.

- Joswiak, Mark, et al. "Dimensionality reduction for visualizing industrial chemical process data." *Control Engineering Practice* 93 (2019): 104189.