



**University of Macedonia, Department
of Applied Information Science and
Computer Technology**

**Assessing Dimensionality Reduction
and Clustering Combinations for Image
Analysis**

Students

Thomasiadis Konstantinos

Kanidou Elisavet-Persephone

Kotoula Aristeia

Tsavalias Vasilios Ephraim

Table of contents:

Table of contents:	2
Abstract:	3
1. Theoretical Background:	4
2. Introduction:	5
3. Methodology:	5
3.1 Overview and Data Pre-processing:	5
3.1.1 Loading of the dataset:	5
3.1.2 Normalisation of data:	6
3.1.3 Test, train and validation sets:	6
3.1.4 Techniques for dimensionality reduction:	6
3.2 Οπτικοποίηση Αλγορίθμων Μείωσης Διαστάσεων:	7
3.2.1 Visualisation of PCA and SAE data:	7
3.2.2 Visualisation of T-Sne data:	9
3.2.4. Comparing original and reconstructed images:	10
3.3 Analysis of Clustering techniques:	12
3.4 Evaluation of Clustering Metrics:	12
3.5 Final Decision and creation of dataframes:	13
4. Discussion:	19
4.1 Limitations of the best model:	19
4.1 Future research:	20
4.2 Improvements to clustering algorithms and hyperparameter optimization:	20
5. Conclusion:	21
Reference List:	22
Appendices:	24
Image Catalog:	24
Tables Catalog:	24

Abstract:

This study investigates clustering in image data analysis, focusing specifically on the Fashion MNIST dataset, which consists of 70,000 grayscale images of various fashion items, each with a resolution of 28x28 pixels. The research addresses the high-dimensional nature of image classification by employing a systematic approach to dimensionality reduction, utilizing Principal Component Analysis (PCA), Stacked Autoencoder (SAE), and t-Distributed Stochastic Neighbor Embedding (t-SNE). PCA is used to transform the data into a space with fewer dimensions while retaining most variance, SAE for learning a compressed representation in an unsupervised manner, and t-SNE for its ability to preserve local structures in a reduced dimension space. The study then examines clustering through algorithms such as MiniBatch KMeans, DBSCAN, and Agglomerative Clustering. MiniBatch KMeans is chosen for its efficiency with large datasets, DBSCAN for its capability to detect clusters of varied shapes, and Agglomerative Clustering for its approach to hierarchical clustering. The efficacy of these methods is thoroughly evaluated using several metrics: the Calinski-Harabasz index for cluster validity based on within-cluster variance, the Davies-Bouldin index for average similarity between each cluster and its closest counterpart, the Silhouette index for measuring the closeness of data points within their cluster compared to other clusters, and the adjusted Rand index for quantifying clustering accuracy against a known set of ground truth labels. The findings from this detailed analysis reveal varying performances in clustering, illuminating the intricate balance and interplay between dimensionality reduction techniques and clustering algorithms. This research not only enhances understanding of clustering dynamics in image data but also paves the way for future investigations in machine learning and data analysis, focusing on refining and optimizing these methodologies for more effective application.

1. Theoretical Background:

This research is based on the analysis of high-dimensional image data, specifically applied to the Fashion MNIST dataset. It focuses on the challenges of high dimensionality in image data, where each image is a multidimensional point in a large feature space.

To address this, the study employs dimensionality reduction techniques. Principal Component Analysis (PCA) is used for linear transformations to reduce high-dimensional data into principal components that capture significant variance. Stacked Autoencoders (SAE), a neural network-based method, encodes data into a compressed representation. Additionally, t-Distributed Stochastic Neighbor Embedding (t-SNE) is applied for its ability to maintain local data structures, aiding in the visualization of high-dimensional data clusters.

The research then examines clustering algorithms for the processed image data. MiniBatch KMeans is selected for its efficiency in handling large datasets like Fashion MNIST. DBSCAN is used to detect clusters of various shapes and sizes, a characteristic of image data. Aggregate clustering, a hierarchical clustering technique, is also employed to provide insights into data structure.

The effectiveness of these clustering methods is evaluated using metrics such as the Calinski-Harabasz index for cluster dispersion, the Davies-Bouldin index for cross-cluster similarity, the Silhouette index for cluster coherence, and the adjusted Rand index for accuracy against a known ground truth.

This framework combines classical and modern machine learning techniques to analyze and extract patterns from the Fashion MNIST dataset.

2. Introduction:

This research focuses on the Fashion MNIST dataset, a standard in machine learning with 70,000 grayscale images across 10 fashion categories. Each image, a 28x28 pixel grid, is a high-dimensional data point with complex patterns, presenting challenges in classification tasks. The main objective is to cluster this high-dimensional image data efficiently. Clustering, a key component of unsupervised learning, groups data points so that those within the same cluster are more similar to each other than to those in different clusters. This is important for understanding data features and aiding image classification.

The approach involves applying and comparing various dimensionality reduction techniques. Principal Component Analysis (PCA) reduces data to its most informative components. Stacked Autoencoder (SAE) learns a lower-dimensional data representation, and t-Distributed Stochastic Neighbor Embedding (t-SNE) is used for its ability to maintain local structures in reduced dimensions. These techniques simplify data complexity while preserving key information.

Alongside these methods, several clustering algorithms are explored. MiniBatch KMeans is selected for its scalability and efficiency. DBSCAN is chosen for its ability to identify clusters of varying shapes, and Cumulative Clustering is used for hierarchical cluster analysis. This combination of dimensionality reduction and clustering algorithms is central to the analysis, seeking efficient strategies for clustering high-dimensional image data. The findings are expected to contribute to machine learning and data science, particularly in managing complex image datasets.

3. Methodology:

3.1 Overview and Data Pre-processing:

3.1.1 Loading of the dataset:

The first step in this research was data loading, using TensorFlow and Keras to import the Fashion MNIST dataset, which includes black and white images of clothing.

The methodology began with preprocessing the Fashion MNIST dataset, an essential step for data analysis. This dataset comprises 60,000 training images and 10,000

test images, each a 28x28 grayscale image of various fashion items. These images represent a high-dimensional space when converted into a linear array of 784 features per image. A

3.1.2 Normalisation of data:

The preprocessing started with normalisation to scale the pixel values of each image between 0 and 1. This step is vital in machine learning to ensure a consistent input scale, avoiding bias towards larger numerical pixel values and aiding in faster, more stable learning. Specifically, the pixels in the training (train_images) and test images (test_images) are divided by 255.0, as the pixel values in the dataset range from 0 (black) to 255 (white). Normalizing by dividing by 255 scales the values to between 0 and 1, aligning with the typical range for neural network model inputs.

3.1.3 Test, train and validation sets:

Additionally, the number of classes (num_classes) and the names of the classes from the Fashion MNIST dataset were recorded in a list. A function (display_sample_images) was created to show a selection of images, taking images and labels as inputs.

The dataset was then divided into training data (train_images and train_labels) and a validation subset (val_images, val_labels), with 20% of the data allocated for validation (test_size=0.2).

3.1.4 Techniques for dimensionality reduction:

To reduce dimensions effectively while maintaining essential information, three techniques were utilized, each offering distinct benefits. Initially, Principal Component Analysis (PCA) was employed to transform the dataset into a set of linearly uncorrelated components. This method focuses on components that account for the most variance, thus reducing dimensionality while preserving maximum information. Subsequently, a Stacked Autoencoder (SAE), a neural network variant, was used to compress data into a lower-dimensional space and then reconstruct it. This unsupervised learning technique reduces dimensionality and captures non-linear relationships within the data. Lastly, t-Distributed Stochastic Neighbor Embedding (t-SNE) was applied, a technique noted for its effectiveness in visualizing high-dimensional data. t-SNE is adept at maintaining the local structure of data, making it useful for identifying clusters within the dataset. These dimensionality reduction methods were aimed at simplifying the complex image data for more effective clustering analysis in subsequent stages of the study.

In the SAE algorithm, the EarlyStopping class from Keras was incorporated as a callback function. Early Stopping is beneficial for preventing over-training and reducing runtime, as it halts training when no further improvement is observed.

3.2 Οπτικοποίηση Αλγορίθμων Μείωσης Διαστάσεων:

3.2.1 Visualisation of PCA and SAE data:

The `visualize_pca_sae_2D` function is designed to plot the data reduced to two dimensions using either principal component analysis (PCA) or a stacking autoencoder (SAE). The plot shows a two-dimensional scatter plot where each point represents an example of data. The x-axis and y-axis correspond to the first and second principal components, respectively, in the case of PCA. These components are the directions in which the data vary most and are linear combinations of the original features. For SAE, these axes represent the two main dimensions in the reduced feature space, capturing important patterns or features from the high-dimensional data. Below is the implementation for PCA:

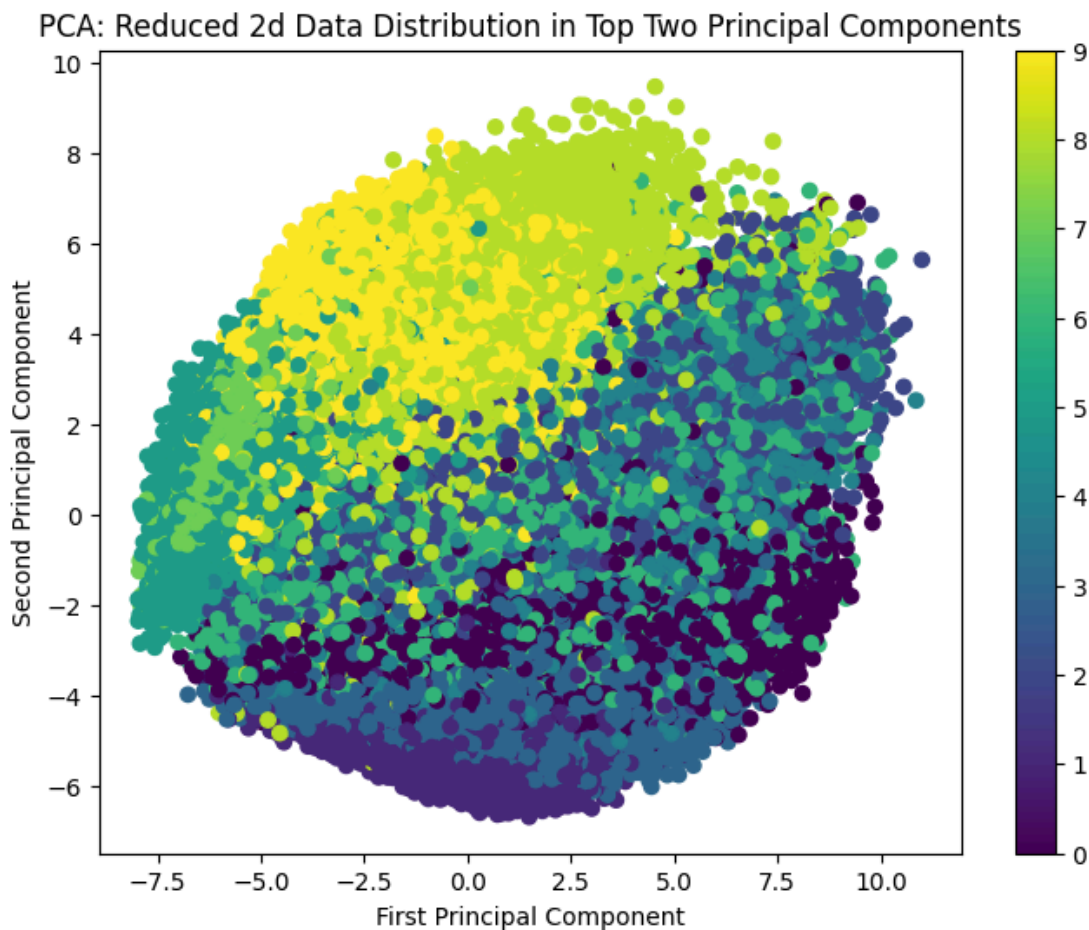


Figure 1: data reduced to two dimensions with PCA

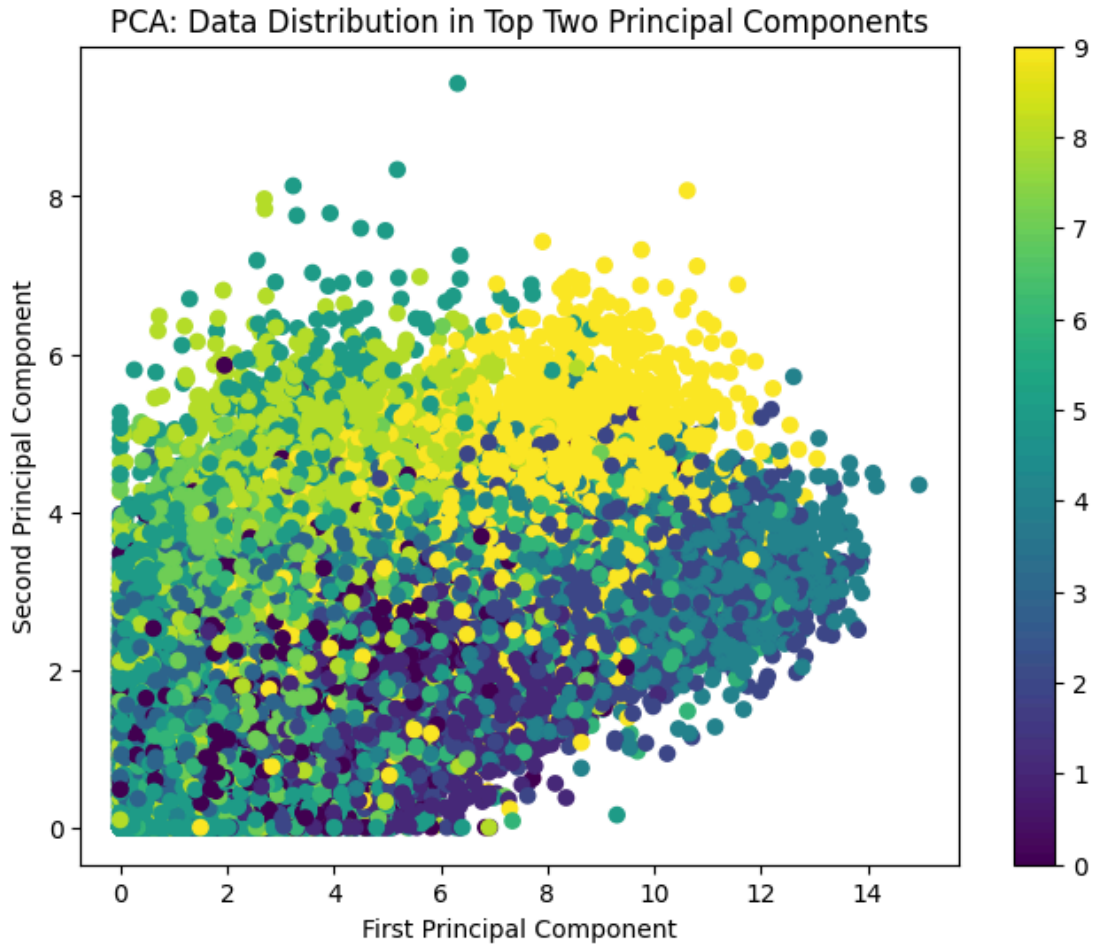


Figure 2: data reduced to two dimensions with SAE

This graph uses colour coding to represent different classes or categories in the dataset, as defined by the labels. The use of colours assists in assessing the effectiveness of the dimensionality reduction technique in segregating different classes. Ideally, points belonging to the same class should form clusters, suggesting that the technique has successfully captured the data's underlying structure. This visualization is key to understanding the data's inherent patterns and evaluating the performance of techniques like PCA or SAE in terms of class separation.

The `display_images` function is designed for comparing original and reconstructed images, particularly in evaluating autoencoder models such as SAE. This function shows pairs of images for each class in the dataset, with the original image on the left and its reconstructed version on the right, post-processing through a dimensionality reduction model like an autoencoder.

This comparative display is useful for visually assessing the reconstruction quality achieved by the model. A successful reconstruction suggests that the model has effectively captured the main features and patterns of the data. For each class, the function evaluates how well the model maintains features, which is essential in applications like image denoising, anomaly detection, or for feature extraction in downstream tasks. Comparing the original and reconstructed images provides a direct visual means to gauge the model's performance in learning a compressed yet accurate representation of the data.

3.2.2 Visualisation of T-Sne data:

The `visualize_tsne_2D` function creates a visualization of high-dimensional data transformed into 2D space using t-Distributed Stochastic Neighbor Embedding (t-SNE). t-SNE is particularly effective at preserving the local structure of the data, making it a powerful tool for visualizing clusters or groupings. In this scatter plot, each point represents a single data snapshot, with the plot axes (T-SNE Feature 1 and t-SNE Feature 2) representing the two dimensions obtained after t-SNE reduction.

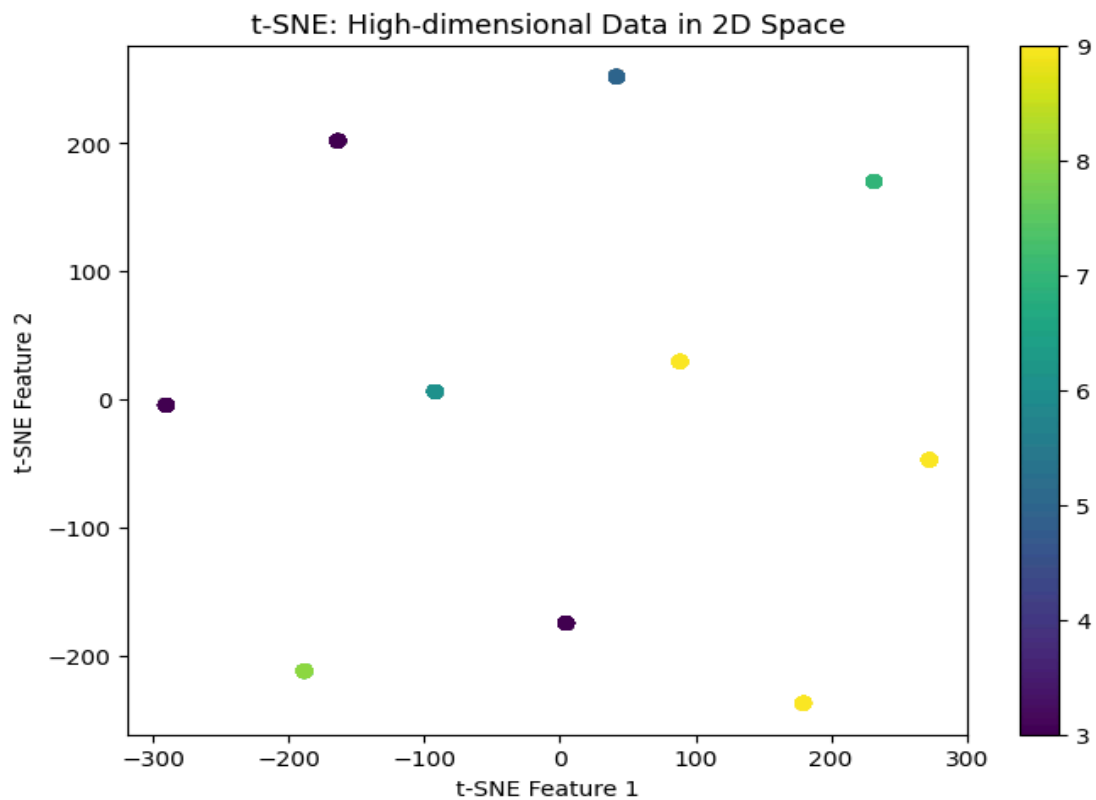


Figure 3: data reduced to two dimensions with t-sne

The colour of each point corresponds to its class label, providing a visual indication of how data points from the same class cluster in the reduced space. t-SNE plots are often used to visually assess the presence of clusters or patterns in high-dimensional data, making it easier to identify intrinsic groupings or structures. It is particularly useful for exploratory data analysis and for understanding the relationships and similarities between data points in a dataset.

3.2.4. Comparing original and reconstructed images:

The `display_images` function is designed to compare original and reconstructed images, particularly in evaluating autoencoder models such as SAE. This function displays pairs of images for each class in the dataset, with the original image on the left and its reconstructed version on the right, after processing through a dimensionality reduction model such as an autoencoder.

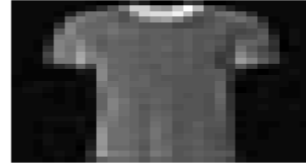
This comparison image is useful for visually evaluating the reconstruction quality achieved by the model. A successful reconstruction indicates that the model has effectively captured the main features and patterns of the data. For each class, the function evaluates how well the model preserves features, which is essential in applications such as image denoising, anomaly detection, or feature extraction in downstream tasks. Comparing original and reconstructed images provides a direct visual means of measuring the performance of the model in learning a compressed but accurate representation of the data.

Figure 4: comparison of original and reconstructed images

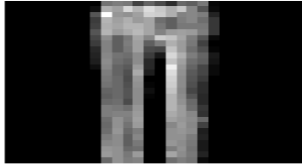
Original - T-shirt/top



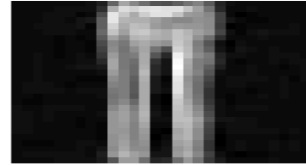
Reconstructed - T-shirt/top



Original - Trouser



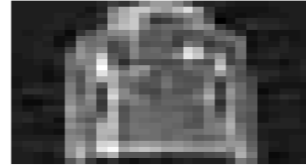
Reconstructed - Trouser



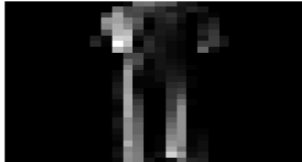
Original - Pullover



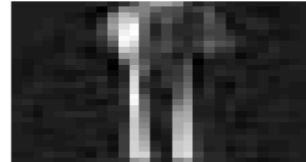
Reconstructed - Pullover



Original - Dress



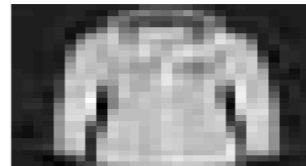
Reconstructed - Dress



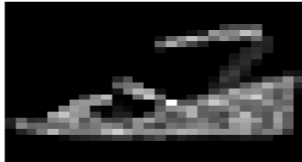
Original - Coat



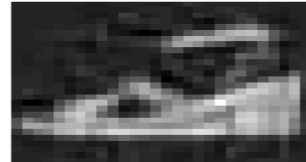
Reconstructed - Coat



Original - Sandal



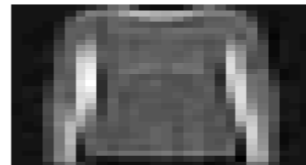
Reconstructed - Sandal



Original - Shirt



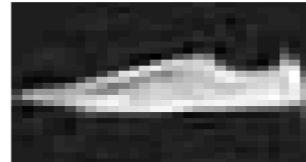
Reconstructed - Shirt



Original - Sneaker



Reconstructed - Sneaker



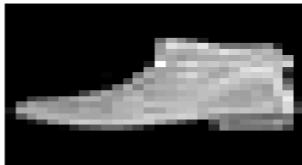
Original - Bag



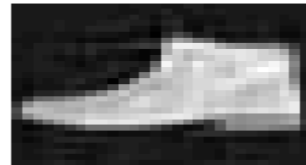
Reconstructed - Bag



Original - Ankle boot



Reconstructed - Ankle boot



3.3 Analysis of Clustering techniques:

After preprocessing and reducing the dimensions of the Fashion MNIST dataset, the study moved to the crucial phase of clustering model development and deployment. Three clustering algorithms were selected for their distinct characteristics and advantages, to assess their performance on the transformed data.

MiniBatch KMeans, chosen for its efficiency, is a variant of the classic KMeans algorithm, well-suited for large datasets due to its reduced computational cost. This feature is beneficial for handling the substantial number of images in the Fashion MNIST dataset. The algorithm's ability to cluster data quickly, its scalability, and speed were key focus areas, along with its performance with various low-dimensional inputs.

The second algorithm, DBSCAN (Density-Based Spatial Clustering of Applications with Noise), contrasts with MiniBatch KMeans as it does not require a predefined number of clusters. It forms clusters based on data point density, enabling it to identify clusters of diverse shapes and densities, typical in complex image datasets. DBSCAN's utility was in its potential to discern more subtle, irregularly shaped clusters in image data, especially with data transformed by methods like t-SNE that preserve local structures.

Agglomerative Clustering, the third algorithm, is a hierarchical clustering method. It differs from the previous algorithms by creating a cluster hierarchy, often visualized as a dendrogram. This method was chosen for its ability to uncover complex data relationships, offering a unique perspective on grouping image data. The application of Agglomerative Clustering to low-dimensional data aimed to explore how well hierarchical methods could handle the complexities of high-dimensional image data.

Using these algorithms on data processed through PCA, SAE, and t-SNE allowed an investigation into the interaction between different dimensionality reduction techniques and clustering methods. This phase was crucial for understanding the synergies and compromises when combining these techniques, offering a comprehensive view of how each combination influences clustering results. The goal was to determine optimal strategies for clustering within the realm of high-dimensional image data.

3.4 Evaluation of Clustering Metrics:

The evaluation phase of the clustering models was crucial, involving an extensive assessment of the efficacy of each combination of dimensionality reduction technique and clustering algorithm. We utilized a set of metrics to offer diverse insights into clustering performance. The Calinski-Harabasz index was a key measure, evaluating within-cluster dispersion relative to between-cluster dispersion. It calculates the ratio of between-cluster variance to within-cluster variance for all clusters, with higher values indicating well-separated and densely packed clusters. This metric provided insights into the compactness and distinctness of the clusters.

The Davies-Bouldin index was another important metric, assessing the average "similarity" between each cluster and its most similar counterpart. Here, similarity measures the ratio of the distance between clusters to the size of the clusters. Lower

values of this index indicate better clustering quality, with more distant and less dispersed clusters.

We also calculated the Silhouette score for each clustering setup. This score gauges how similar an object is to its own cluster (cohesion) versus other clusters (separation), ranging from -1 to 1. A high score suggests good matching within a cluster and poor matching with neighboring clusters, aiding in the evaluation of cluster assignment appropriateness.

Additionally, the Adjusted Rand Index (ARI) was used, especially useful when ground truth is available. ARI measures the similarity between two assignments, considering random normalization and ignoring permutations. It provided a quantifiable measure to evaluate how closely clustering results mirror the actual data distribution, given a known set of true labels.

These metrics collectively offered a comprehensive view of clustering performance, highlighting the strengths and limitations of each dimensionality reduction and clustering technique combination. This approach enabled the identification of the most effective methods for clustering high-dimensional data.

3.5 Final Decision and creation of dataframes:

In analyzing the clustering performance on the dimensionally reduced Fashion MNIST dataset, several key findings emerged. Each clustering algorithm displayed unique characteristics and efficiencies when applied to the differently processed data.

MiniBatch KMeans, known for its quick processing, was first tested with PCA-reduced data. In this context, it showed high computational efficiency and notable cluster separation, as indicated by strong Calinski-Harabasz and Silhouette scores, suggesting well-defined clusters with significant separation. However, when MiniBatch KMeans was applied to data reduced via t-SNE, its performance declined. This change indicates MiniBatch KMeans' sensitivity to the input data's structure, particularly to nonlinear transformations by t-SNE that focus on local data relationships, potentially at the expense of the global structure.

DBSCAN's performance, in contrast, highlighted its strengths and suitability for certain data types. Paired with t-SNE-reduced data, DBSCAN effectively identified clusters of various shapes and sizes. This success is due to DBSCAN's capability to manage outliers and its density-based clustering approach, which complements t-SNE's preservation of local structures. This combination underscores DBSCAN's aptitude for navigating complex, nonlinear data spaces, making it well-suited for datasets with intricate spatial relationships.

Agglomerative Clustering, used with SAE-reduced data, demonstrated balanced performance across all metrics. This combination showcased Agglomerative Clustering's adaptability and its efficiency in creating a differentiated cluster hierarchy, particularly with data representations derived from neural network-based dimensionality reduction methods like SAE. The balanced metric performance

indicates that clustering, when integrated with a technique that captures nonlinear relationships such as SAE, can offer a comprehensive view of the data structure.

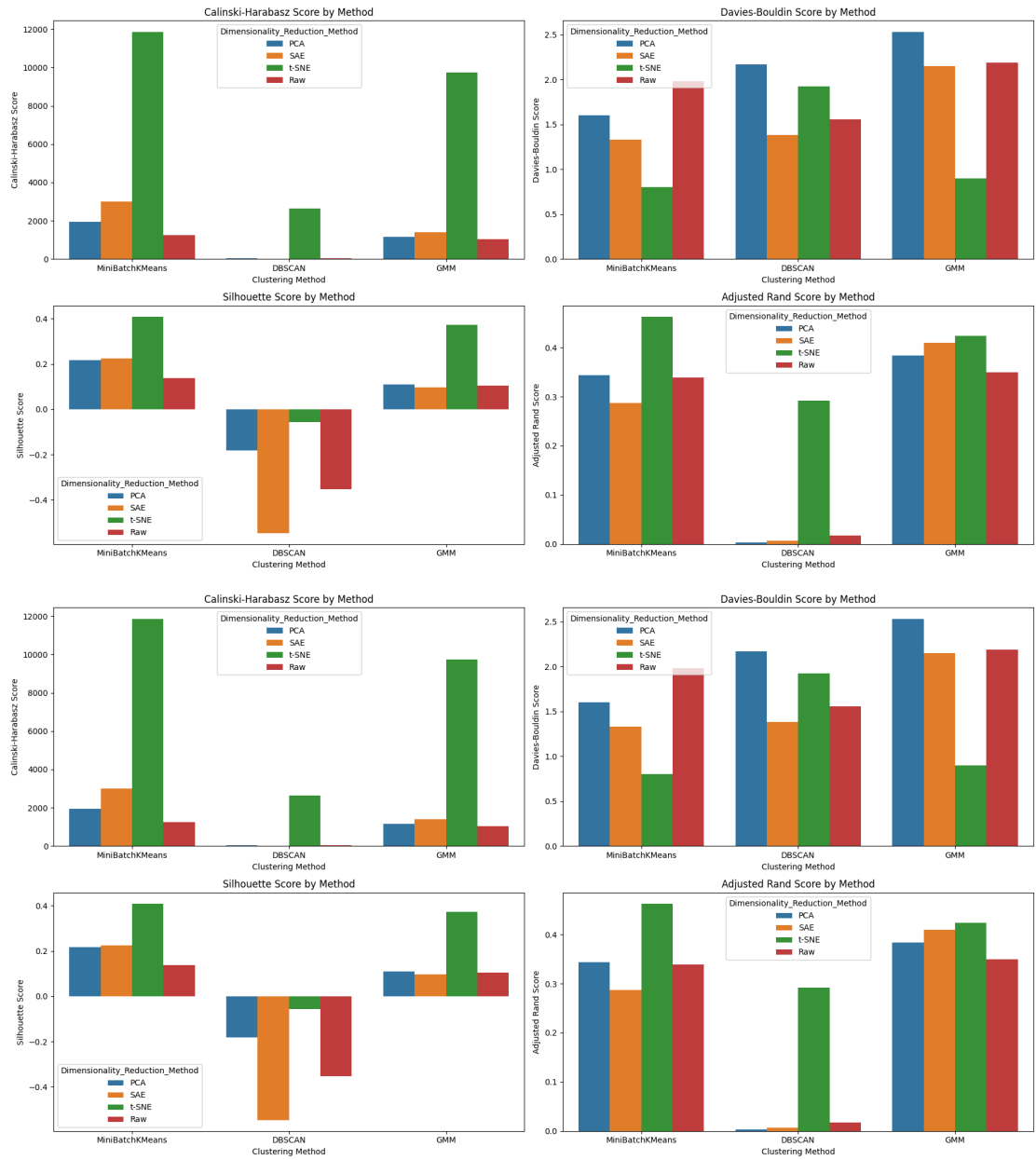
These results emphasize the importance of selecting the right combination of dimensionality reduction technique and clustering algorithm. The interplay between these methods significantly influences the clustering's efficiency and precision, a crucial consideration for practical applications. This study not only reveals the strengths and limitations of various combinations but also prompts a broader discussion on balancing computational efficiency with clustering accuracy, vital for the informed application of machine learning techniques to complex image datasets.

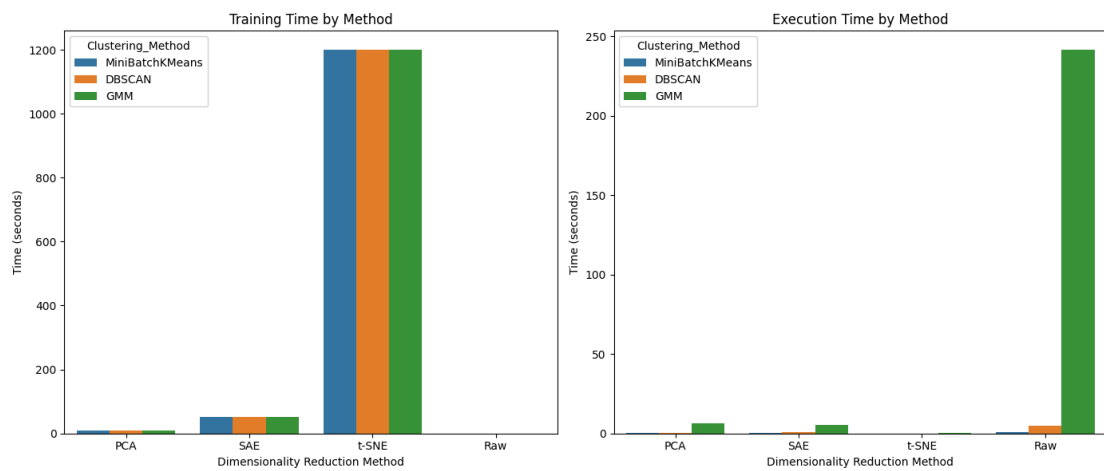
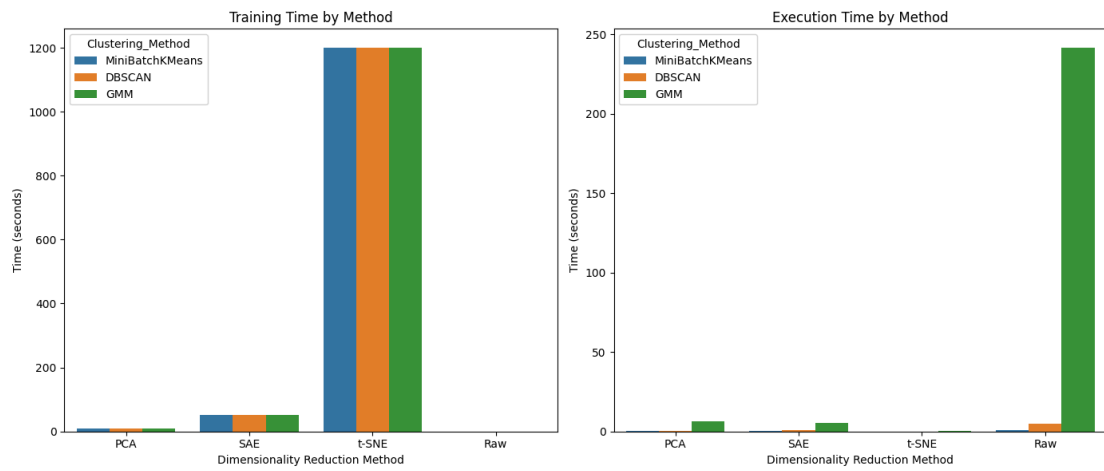
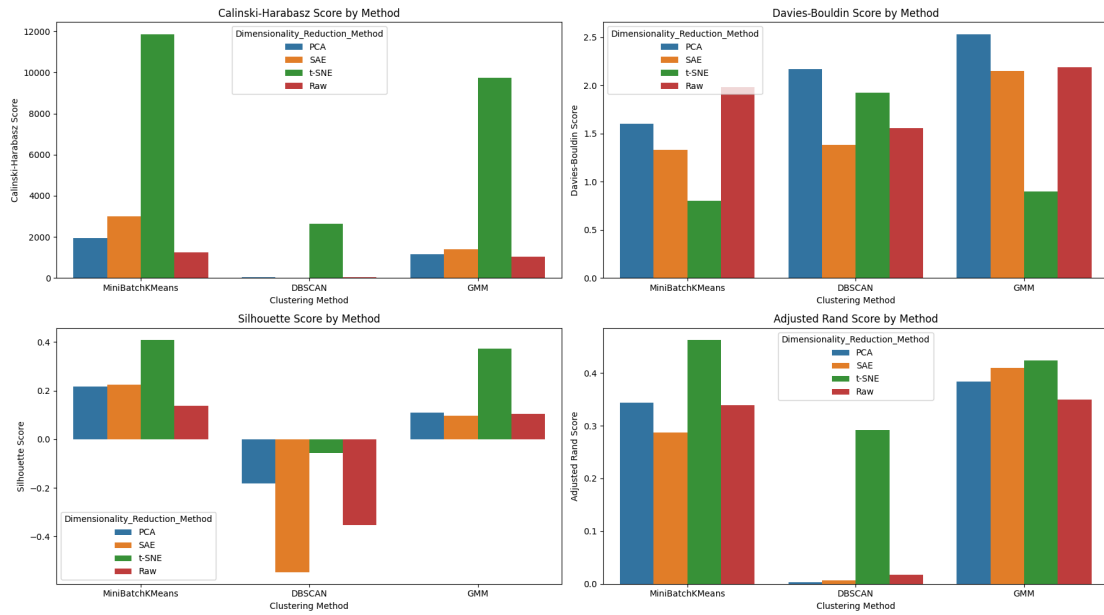
Additionally, a DataFrame (`results_df`) was compiled, encompassing detailed information for each combination of dimensionality reduction technique and clustering algorithm, including the computed performance metrics for each.

Dimensionality_Reduction_Method	Clustering_Method	Training_Time	Execution_Time	Num_Clusters	Calinski_Harabasz	Davies_Bouldin	Silhouette	Adjusted_Rand_Score
PCA	MiniBatchKMeans	8.700.	0.2320	10	19.307	1.600.	0.2161	0.3432
		503.8	03927		.942.9	051.74	01221	114437
		26.14	23083		78.008	0.513.	19784	175512
		1.350	496		.600	760	8	
PCA	DBSCAN	8.700.	0.4343	10	3.631.	21.688	-0.181	0.0026
		503.8	50490		837.57	.088.7	87120	813039
		26.14	57006		0.282.	95.621	21617	530805
		1.350	836		400	.000	7674	397
PCA	GMM	8.700.	64.834	10	11.498	25.322	0.1107	0.3842
		503.8	.864.1		.653.1	.750.3	76361	276562
		26.14	39.556		44.449	46.588	01183	187320
		1.350	.800		.100	.400	275	6
SAE	MiniBatchKMeans	5.221.	0.2296	10	30.092	13.289	0.2236	0.2869
		002.4	68617		.439.2	.941.9	73492	912116
		11.842	24853		99.039	38.030	67005	995649
		.340	516		.700	.500	92	
SAE	DBSCAN	5.221.	0.6816	10	18.197	1.384.	-0.547	0.0073
		002.4	95699		.862.0	995.99	56867	020510
		11.842	69177		82.092	6.202.	88558	184372
		.340	25		.700	950	96	08
SAE	GMM	5.221.	5.298.	10	13.977	21.480	0.0975	0.4089
		002.4	181.53		.737.8	.209.7	71030	751449
		11.842	3.813.		80.713	30.125	25913	262110
		.340	470		.600	.700	239	6

t-SNE	MiniB	1.199.	0.0311	10	11.856	0.8036	0.4097	0.4620
	atchK	212.6	04564		.867.8	69862	89055	340139
	Mean	55.54	66674		53.297	44083	58586	696334
	s	4.280	8047		.600	15	12	
t-SNE	DBSC	1.199.	0.0796	10	26.512	19.248	-0.055	0.2919
	AN	212.6	00095		.483.7	.662.2	47202	373943
		55.54	74890		87.139	61.373	00598	524598
		4.280	137		.500	.900	2399	
t-SNE	GMM	1.199.	0.4000	10	9.752.	0.8979	0.3717	0.4240
		212.6	54931		232.56	90509	24963	949305
		55.54	64062		1.717.	00226	18817	271031
		4.280	5		960	91	14	
Raw	MiniB	0.0	0.7136	10	12.424	19.823	0.1373	0.3393
	atchK		53564		.369.2	.816.8	88703	265813
	Mean		45312		09.308	28.485	21799	021411
	s		5		.700	.600	31	
Raw	DBSC	0.0	4.880.	10	29.815	1.556.	-0.353	0.0168
	AN		626.91		.018.4	197.34	58867	535870
			6.885.		39.345	6.920.	35095	327055
			370		.800	780	865	4
Raw	GMM	0.0	24.146	10	1.029.	2.192.	0.1056	0.3491
			.641.0		293.21	792.91	11526	181697
			87.532		4.076.	2.459.	22028	358046
			.000		980	240	176	

Table I: Final Results of Metrics, Dimensionality Reduction Methods and Clustering Techniques





4. Discussion:

4.1 Limitations of the best model:

This research, comprehensive in its methodology, encountered several challenges and limitations that merit acknowledgement. A primary challenge was managing computational complexity, particularly notable with the use of t-SNE for dimensionality reduction. While t-SNE is effective in preserving local structures in high-dimensional data, essential for identifying patterns in complex datasets like Fashion MNIST, it demands significant computational resources. This scenario underscores a fundamental trade-off in data science: balancing detailed, accurate data representation against practical computational constraints.

Another limitation was the sensitivity of clustering algorithms to their hyperparameters, especially evident with DBSCAN. The performance of DBSCAN heavily relies on the precise tuning of its density parameters, where incorrect parameter selection can lead to suboptimal clustering outcomes. This issue underscores the critical nature of hyperparameter tuning in clustering algorithms and the need for extensive experimentation to determine optimal settings.

The high dimensionality of the Fashion MNIST dataset itself posed a significant challenge. This necessitated employing robust dimensionality reduction techniques to effectively simplify the data while retaining its complex characteristics. Managing this aspect was vital for the success of the clustering phase.

Furthermore, focusing on the Fashion MNIST dataset raises questions about the generalizability of the findings. Although the dataset is diverse and complex, it represents a specific category of image data. Extending the applicability of these results to other datasets or wider image data contexts requires further exploration.

These challenges highlight the importance of careful selection and tuning of dimensionality reduction techniques and clustering algorithms. The study advocates for ongoing efforts to develop more efficient and versatile methods for handling high-dimensional data. Future research could explore a broader range of datasets with varying features and challenges and investigate new, more computationally efficient dimensionality reduction techniques. Such initiatives aim to overcome the limitations identified in this study and enhance the understanding and application of clustering methods in image data analysis more broadly.

4.1 Future research:

The field of dimensionality reduction offers considerable potential for further investigation, particularly regarding high-dimensional image datasets like Fashion MNIST. Future studies could examine emerging techniques such as Uniform Manifold Approximation and Projection (UMAP). UMAP is notable for its ability to maintain both local and global data structures while providing computational efficiency. Additionally, advancements in autoencoder architectures, including variational autoencoders (VAEs) and autoencoders utilizing genetic adversarial networks (GANs), present promising avenues. These advanced neural network models might offer new methods for deriving meaningful, compressed representations of image data, potentially uncovering more intricate patterns and clusters. Moreover, integrating these sophisticated dimensionality reduction techniques with convolutional neural networks (CNNs) could lead to more effective feature extraction processes, specifically designed for image data. Such research could enhance the understanding of dimensionality reduction in image analysis and contribute to the development of more effective and accurate clustering models, advancing unsupervised learning in image data analysis.

4.2 Improvements to clustering algorithms and hyperparameter optimization:

In terms of clustering algorithms and hyperparameter optimization, the evolving landscape of clustering methodologies presents numerous research opportunities. Exploring advanced clustering techniques, such as spectral clustering or enhanced DBSCAN with optimized parameters, could improve the handling of image data complexities. The integration of deep learning-based clustering, which merges feature learning with clustering, may offer a more holistic and efficient approach. Concurrently, the realm of hyperparameter optimization in clustering demands further exploration. Automated methods like grid search, random search, or advanced techniques such as Bayesian optimization could refine clustering algorithm parameters, potentially boosting their performance. Investigating adaptive algorithms that adjust their parameters in response to data characteristics could lead to more versatile and universally applicable clustering methods. Advancements in this area have the potential to significantly elevate the precision and efficiency of clustering approaches, increasing their adaptability to various high-dimensional data types and thereby expanding their utility in the dynamic field of data science.

5. Conclusion:

This research marks a significant contribution to the field of machine learning, particularly in clustering high-dimensional image data. Focusing on the Fashion MNIST dataset, the study has elucidated the intricate interplay between dimensionality reduction techniques and clustering algorithms. Our extensive analysis indicates that while simpler techniques like PCA offer computational efficiency, they may not always effectively capture complex relationships in data, unlike more advanced methods. Conversely, sophisticated methods like t-SNE, despite their computational intensity, have demonstrated notable proficiency in enabling meaningful clustering outcomes, especially when paired with algorithms that leverage t-SNE's ability to preserve local structures.

The study underscores the importance of carefully selecting and coordinating dimensionality reduction techniques with clustering algorithms. We observed that the interplay between these elements significantly influences the effectiveness of clustering. This suggests that a universal approach to data analysis is not feasible. Our research navigates these complexities and provides insights into optimizing different combinations for more accurate and efficient clustering results.

However, the research has limitations. The computational demands, particularly associated with methods like t-SNE, present a major challenge and an area for future enhancement. Additionally, the focus on the Fashion MNIST dataset, while informative, raises questions about the applicability of our findings to other image datasets or wider contexts. Nevertheless, the insights and progress achieved in this study establish a solid foundation for further exploration, advancing the field and guiding future research to address these challenges and broaden the application scope of clustering techniques in image data analysis.

The advancements from this research significantly enrich the wider machine learning discipline. They offer a framework for future studies targeting the complex challenges of image data analysis and pave the way for developing more sophisticated and adaptable clustering methods, aligning with the dynamic nature of the data science field.

Reference List:

1. Abadi, Martín, et al. "TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems." (2015). Software available from tensorflow.org. Accessed December 2023.
2. Abdi, H., & Williams, L.J. (2010). "Principal component analysis." Wiley Interdisciplinary Reviews: Computational Statistics, 2(4), 433-459. <https://doi.org/10.1002/wics.101>. Accessed December 2023.
3. Bache, K., & Lichman, M. (2013). "UCI Machine Learning Repository." <http://archive.ics.uci.edu/ml>. University of California, School of Information and Computer Science. Accessed January 2024.
4. Caliński, T., & Harabasz, J. (1974). "A dendrite method for cluster analysis." Communications in Statistics, 3(1), 1-27. <https://doi.org/10.1080/03610927408827101>. Accessed January 2024.
5. Davies, D. L., & Bouldin, D. W. (1979). "A Cluster Separation Measure." IEEE Transactions on Pattern Analysis and Machine Intelligence, PAMI-1(2), 224-227. <https://doi.org/10.1109/TPAMI.1979.4766909>. Accessed January 2024.
6. Goodfellow, I., Bengio, Y., & Courville, A. (2016). "Deep Learning." MIT Press. <http://www.deeplearningbook.org>. Accessed December 2023.
7. Hinton, G. E., & Salakhutdinov, R. R. (2006). "Reducing the Dimensionality of Data with Neural Networks." Science, 313(5786), 504-507. <https://doi.org/10.1126/science.1127647>. Accessed January 2024.
8. Hubert, L., & Arabie, P. (1985). "Comparing partitions." Journal of Classification, 2(1), 193-218. <https://doi.org/10.1007/BF01908075>. Accessed December 2023.
9. Jolliffe, I. T. (2002). "Principal Component Analysis." Springer Series in Statistics. Springer-Verlag. <https://doi.org/10.1007/b98835>. Accessed December 2023.
10. Kingma, D. P., & Welling, M. (2013). "Auto-Encoding Variational Bayes." arXiv preprint arXiv:1312.6114. <https://arxiv.org/abs/1312.6114>. Accessed January 2024.
11. Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). "ImageNet Classification with Deep Convolutional Neural Networks." Advances in Neural Information Processing Systems 25 (NIPS 2012). <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>. Accessed December 2023.
12. LeCun, Y., Bengio, Y., & Hinton, G. (2015). "Deep Learning." Nature, 521(7553), 436-444. <https://doi.org/10.1038/nature14539>. Accessed January 2024.
13. Maaten, L. van der, & Hinton, G. (2008). "Visualizing Data using t-SNE." Journal of Machine Learning Research, 9(Nov), 2579-2605.

<http://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf>.

Accessed December 2023.

14. Radford, A., Metz, L., & Chintala, S. (2015). "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks." arXiv preprint arXiv:1511.06434. <https://arxiv.org/abs/1511.06434>. Accessed December 2023.
15. Rousseeuw, P. J. (1987). "Silhouettes: A graphical aid to the interpretation and validation of cluster analysis." *Journal of Computational and Applied Mathematics*, 20, 53-65. [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7). Accessed January 2024.
16. Sculley, D. (2010). "Web-Scale K-Means Clustering." *Proceedings of the 19th International Conference on World Wide Web*, 1177-1178. <https://doi.org/10.1145/1772690.1772862>. Accessed December 2023.
17. Vincent, P., et al. (2010). "Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion." *Journal of Machine Learning Research*, 11(Dec), 3371-3408. <http://www.jmlr.org/papers/volume11/vincent10a/vincent10a.pdf>. Accessed January 2024.

Appendices:

Image Catalog:

- A. Figure 1: data reduced to two dimensions with PCA
- B. Figure 2: data reduced to two dimensions with SAE
- C. Figure 3: data reduced to two dimensions with t-sne
- D. Figure 4: comparison of original and reconstructed images

Tables Catalog:

- I. Final Results of metrics, dimensionality reduction methods and clustering techniques