

# Assignment 2 Report

## Introduction

The MNIST dataset, a benchmark in the field of machine learning and computer vision, has long served as a foundational resource for developing and testing various algorithms. Comprising 70,000 28x28 pixel images of handwritten digits (0 through 9), MNIST has become synonymous with image classification tasks. In this report, we delve into the exploration of dimensionality reduction techniques, specifically focusing on Principal Component Analysis (PCA), and extend our analysis to clustering algorithms such as K-means and Gaussian Mixture Model (GMM).

## Objectives:

- Investigate the effectiveness of K-means clustering on the raw MNIST dataset, providing insights into its limitations and strengths.
- Develop a deeper understanding of PCA, implement it from scratch using Singular Value Decomposition (SVD), and identify the optimal number of components for dimensionality reduction.
- Apply GMM clustering to the PCA-transformed data, examining the resulting clusters in a lower-dimensional space.
- Compare and contrast the results obtained through these methodologies, shedding light on the impact of dimensionality reduction on clustering performance.

Through this exploration, we aim to provide a comprehensive view of the interplay between dimensionality reduction and clustering algorithms, offering valuable insights for practitioners and researchers alike.

## Background

*2.1 MNIST Dataset Overview:* The MNIST dataset, originating from the Modified National Institute of Standards and Technology, stands as a cornerstone in the machine learning community. Comprising 70,000 grayscale images, each depicting a handwritten digit from 0 to 9, MNIST serves as a quintessential resource for training and evaluating algorithms, particularly those geared towards image classification. With 60,000 images designated for training and 10,000 for testing, MNIST encapsulates the diversity of human handwriting, challenging algorithms to recognize and categorize digits accurately.

*2.2 Importance of Preprocessing:* Before delving into dimensionality reduction and clustering, it's crucial to highlight the significance of preprocessing, specifically data normalization. The pixel values in the MNIST images, originally ranging from 0 to 255, are scaled to the interval  $[0, 1]$ . Normalization ensures that each feature contributes proportionally to the analysis, preventing variables with larger scales from dominating the results. This step lays the foundation for subsequent analyses and aids in achieving reliable and meaningful outcomes.

*2.3 Challenges in Dimensionality:* The MNIST dataset, with its 28x28 pixel images, initially presents itself as a high-dimensional space, where each pixel contributes to a separate dimension. This high dimensionality poses challenges such as increased computational complexity, susceptibility to overfitting, and the curse of dimensionality. As a result, exploring dimensionality reduction techniques becomes imperative, with the goal of retaining essential information while mitigating these challenges.

In the subsequent sections of this report, we will address these challenges and delve into the intricacies of dimensionality reduction using PCA, shedding light on its efficacy in capturing essential patterns within the MNIST dataset.

## **K-Means Clustering on MNIST**

*3.1 Overview of K-Means Clustering:* K-Means clustering is a widely used unsupervised learning algorithm that aims to partition a dataset into K distinct, non-overlapping subgroups or clusters. The algorithm iteratively assigns data points to clusters based on the mean of the feature values, converging to a solution where each point belongs to the cluster with the nearest mean.

*3.2 Application to MNIST:* Applying K-Means clustering to the raw MNIST dataset involves grouping the 28x28 pixel images into clusters based on their feature similarities. Varying the number of clusters (K) allows us to explore different partitionings of the dataset. We have experimented with K values of 10, 7, and 4 to observe how the algorithm categorizes the diverse handwritten digits.

*3.3 Visualizations and Interpretations:* Visualizations of the resulting clusters provide insights into the grouping patterns discovered by K-Means. Each cluster represents a set of handwritten digits that share common characteristics. Interpretations involve analyzing the dominant digit types within each cluster, understanding the distribution of digits across clusters, and noting any challenges or limitations in the algorithm's ability to differentiate between certain digit classes.

*3.4 Comparison with Previous Tasks:* In this section, we compare the outcomes of K-Means clustering on the MNIST dataset with the results obtained from the previous task. It involves contrasting the cluster characteristics, sizes, and interpretability between the Euclidean distance-based clustering and the subsequent dimensionality reduction approaches.

Through the examination of K-Means clustering on MNIST, we aim to uncover insights into the inherent patterns within the dataset and gain a preliminary understanding of its clustering behavior.

## **PCA from Scratch**

*4.1 Explanation of Principal Component Analysis (PCA):* Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms high-dimensional data into a lower-dimensional space while retaining the most significant information. The principal components are linear combinations of the original features, ordered by their ability to explain variance in the data. By selecting a subset of these components, one can achieve effective dimensionality reduction.

*4.2 Implementation of PCA from Scratch:* The implementation of PCA involves Singular Value Decomposition (SVD), a method for factorizing a matrix into three separate matrices. The SVD process is applied to the MNIST dataset to decompose it into the left singular vectors ( $U$ ), singular values ( $S$ ), and right singular vectors ( $V^t$ ). By selecting a subset of these components, we construct the PCA-transformed dataset.

*4.3 Visualization of Explained Variance:* Explained variance ratio and cumulative explained variance are crucial metrics for determining the optimal number of components. Visualizations of these metrics provide insights into how much information is retained as we increase the number of components. A scree plot or elbow plot is generated to help identify the point of diminishing returns, indicating the optimal number of components.

*4.4 Determination of Optimal Components:* The optimal number of components is determined based on a chosen threshold for cumulative explained variance. A threshold of 95% is commonly used, ensuring that the selected components capture a substantial portion of the dataset's variability while avoiding overfitting.

In this section, we lay the foundation for subsequent analyses by implementing PCA from scratch, visualizing its results, and determining the optimal number of components. The insights gained will guide the subsequent application of GMM clustering on the PCA-transformed data.

## **GMM Clustering on PCA-Transformed Data**

*5.1 Application of Gaussian Mixture Model (GMM):* Gaussian Mixture Model (GMM) is a probabilistic model that represents a dataset as a mixture of several Gaussian distributions. Each distribution is associated with a cluster, and the model estimates the parameters such as mean, covariance, and weight for each cluster. Applying GMM to the PCA-transformed data involves capturing the underlying probability distribution and uncovering intricate structures that may not be apparent in the original feature space.

*5.2 Visualizations of GMM Clusters:* Visualizations of GMM clusters in the lower-dimensional space created by PCA provide a nuanced view of the data's inherent structures. Each cluster is represented by a Gaussian distribution, and the visualizations offer insights into the density and shape of the clusters. The choice of the number of clusters ( $K$ ) in GMM affects the granularity and complexity of the discovered patterns.

*5.3 Interpretation of Cluster Characteristics:* Interpreting the characteristics of clusters involves analyzing the digit types that dominate each cluster, understanding the density and spread of data points within clusters, and discerning any discernible patterns. Comparisons with the results from K-Means clustering and observations about the flexibility of GMM in capturing non-linear structures are discussed.

*5.4 Impact of Dimensionality Reduction:* This section addresses the impact of dimensionality reduction on the clustering performance. By comparing the GMM clustering results on PCA-transformed data with the earlier K-Means clustering results, we gain insights into how dimensionality reduction affects the ability to capture complex structures and relationships within the MNIST dataset.

Through the application of GMM clustering on PCA-transformed data, we aim to uncover more nuanced cluster characteristics and assess the effectiveness of the probabilistic modeling approach in capturing underlying data distributions.

## **Comparison of Results**

*6.1 Comparative Analysis of K-Means and GMM:* This section presents a detailed comparison of the results obtained from K-Means clustering and GMM clustering on the MNIST dataset. Key aspects such as cluster characteristics, interpretability, and the impact of dimensionality reduction are thoroughly examined.

*6.2 Clustering Flexibility:* Discussing the flexibility of each algorithm in capturing different types of clusters is essential. While K-Means assumes spherical clusters and uniform variance, GMM, being a probabilistic model, can handle clusters with varying shapes and sizes. This section delves into the implications of these differences on the clustering outcomes.

*6.3 Visual Comparisons:* Visual comparisons, including side-by-side visualizations of K-Means and GMM clusters, aid in understanding the disparities in their grouping strategies. Highlighting specific instances where GMM excels in capturing complex patterns that K-Means might overlook adds depth to the analysis.

*6.4 Interpretability Challenges:* Interpretability is a critical aspect of clustering results. Discussing the interpretability challenges posed by GMM, where clusters are represented by Gaussian distributions, contrasts with the more straightforward interpretability of K-Means clusters.

*6.5 Dimensionality Reduction Trade-offs:* The trade-offs associated with dimensionality reduction through PCA are explored in this section. While PCA aids in computational efficiency and visualization, it also raises questions about information loss and the balance between retaining essential patterns and minimizing noise.

Through this comprehensive comparative analysis, we aim to provide a nuanced understanding of the strengths and limitations of K-Means clustering and GMM clustering, shedding light on their applicability to the MNIST dataset and similar real-world scenarios.

## **Optimal Number of Components and Interpretability Challenges**

*7.1 Optimal Number of Components:* Determining the optimal number of components for PCA is a crucial step in achieving effective dimensionality reduction. In this section, we elaborate on the methodology used to identify the optimal number of components based on the explained variance ratio and cumulative explained variance. The chosen threshold of 95% provides a balance between retaining sufficient information and avoiding overfitting.

*7.2 Insights from Explained Variance Visualizations:* Visualizations of the explained variance ratio and cumulative explained variance offer insights into the distribution of variance across principal components. Analyzing these visualizations aids in understanding how rapidly information is captured as we increase the number of components and identifies the point where additional components contribute marginally to the overall variance.

*7.3 Challenges in Interpretability:* While PCA provides an effective means of dimensionality reduction, the interpretability of the resulting principal components can be challenging. Each principal component is a linear combination of the original features, making it less intuitive to associate specific components with meaningful patterns. This section addresses the interpretability challenges posed by PCA and explores strategies to enhance the understanding of the transformed features.

*7.4 Visualization of PCA-Transformed Data:* Visualizing the PCA-transformed data in lower-dimensional space offers a qualitative understanding of how well the selected components capture the variability in the original dataset. Scatter plots and other visualizations

demonstrate the distribution of data points in the reduced feature space and highlight any discernible patterns.

Through a detailed exploration of the optimal number of components and the interpretability challenges associated with PCA, this section contributes to the foundational understanding of the dimensionality reduction process and its implications for subsequent analyses, including clustering.

## **Limitations, Conclusion, and Future Work**

*8.1 Limitations of the Approach:* Discussing the limitations of the analysis conducted is crucial for a comprehensive understanding. Addressing aspects such as the assumptions made, challenges encountered, and potential sources of bias or error provides context for interpreting the results. Limitations may include sensitivity to hyperparameter choices, the reliance on predefined thresholds, and the assumptions inherent in the clustering algorithms.

*8.2 Reflection on Interpretability:* Reflecting on the interpretability challenges posed by both dimensionality reduction and clustering methodologies, this section delves into the broader implications for real-world applications. It discusses the trade-offs between model complexity and interpretability and emphasizes the need for a nuanced understanding of the limitations associated with unsupervised learning techniques.

*8.3 Conclusion:* Summarizing the key findings and insights derived from the analysis, this section provides a concise conclusion. It revisits the main objectives outlined at the beginning of the report and highlights the contributions of the study to the broader understanding of dimensionality reduction and clustering on the MNIST dataset. The conclusion serves as a synthesis of the knowledge gained and sets the stage for potential applications and future research directions.



*8.4 Future Work:* Identifying avenues for future research is integral to the conclusion. This section outlines potential extensions of the current study, including exploring alternative dimensionality reduction techniques, incorporating additional clustering algorithms, or adapting the methodology for diverse datasets. Addressing the limitations identified in the analysis, this section provides a roadmap for researchers and practitioners interested in advancing the field.

In closing, this page encapsulates the broader context of the study, acknowledges its limitations, and paves the way for future investigations in the dynamic and evolving landscape of dimensionality reduction and clustering.

## **Comparative Evaluation and Insights**

*9.1 Comparative Evaluation Metrics:* To comprehensively assess the performance of both K-Means clustering and GMM clustering on PCA-transformed data, this section introduces and applies relevant evaluation metrics. Metrics such as silhouette score, Davies-Bouldin index, or other suitable measures are employed to quantify the compactness and separation of clusters. A comparative analysis using these metrics offers a quantitative perspective on the effectiveness of the two methodologies.

*9.2 Validation and Robustness:* To ensure the validity and robustness of the clustering results, considerations of sensitivity to initialization, the impact of hyperparameters, and the stability of the clusters are addressed. Robustness checks, sensitivity analyses, or cross-validation approaches may be employed to enhance the reliability of the findings.

*9.3 Domain-Specific Considerations:* Discussing how the results align with domain-specific knowledge, if available, enhances the practical relevance of the analysis. For instance, in the

context of the MNIST dataset, understanding how well the clustering results align with human perceptions of digit similarity provides additional context and validation.

*9.4 Implications for Real-world Applications:* Translating the findings into practical implications for real-world applications is crucial. This section discusses how the insights derived from the clustering analysis could inform decision-making processes, contribute to pattern recognition systems, or be utilized in various domains such as image recognition, fraud detection, or customer segmentation.

Through a comprehensive evaluation and reflection on the insights gained, this section provides a deeper understanding of the practical implications and potential applications of the clustering methodologies applied to the MNIST dataset.

## **Conclusion and Key Takeaways**

*10.1 Recapitulation of Findings:* This section succinctly summarizes the key findings and insights obtained throughout the report. It revisits the main objectives outlined in the introduction and highlights the significant outcomes of the analysis, focusing on both the results of clustering methodologies and the impact of dimensionality reduction.

*10.2 Contribution to the Field:* Reflecting on the broader contribution of the study to the field of machine learning and clustering, this section emphasizes any novel insights, methodologies, or perspectives introduced. It discusses how the analysis extends existing knowledge and contributes to advancing our understanding of dimensionality reduction and clustering on complex datasets like MNIST.

*10.3 Practical Recommendations:* Based on the analysis conducted, practical recommendations are provided for researchers, practitioners, and machine learning enthusiasts. These recommendations may include insights into the choice of clustering algorithms, considerations

for dimensionality reduction, and potential strategies for enhancing the interpretability of clustering results.

*10.4 Final Thoughts:* Concluding the report, this section offers final thoughts on the significance of the study and its implications for future research in machine learning, dimensionality reduction, and clustering. It may include reflections on the challenges encountered, lessons learned, and the broader impact of the study on the ever-evolving landscape of data science.

*10.5 Acknowledgments and References:* Any acknowledgments for contributions, guidance, or support received during the course of the study are included in this section. Additionally, a comprehensive list of references is provided, acknowledging the works and studies that informed the methodologies and concepts discussed in the report.

In closing, this page encapsulates the essence of the report, summarizes its contributions, and provides a roadmap for future research and applications in the dynamic field of machine learning and clustering.