

# **Explain Clustering and Dimensionality Reduction**

**Bachelor of Technology  
Computer Science and Engineering**

Submitted By

ARKAPRATIM GHOSH (13000121058)

**Neural Networks and Deep Learning  
PCC-CS702A**

SEPTEMBER 2024



**Techno Main Salt Lake  
EM-4/1, Sector-V,  
Kolkata- 700091  
West Bengal  
India**

## Table of Contents

1. Introduction.....	3
2. Body.....	3
3. Conclusion.....	7
4. References.....	7

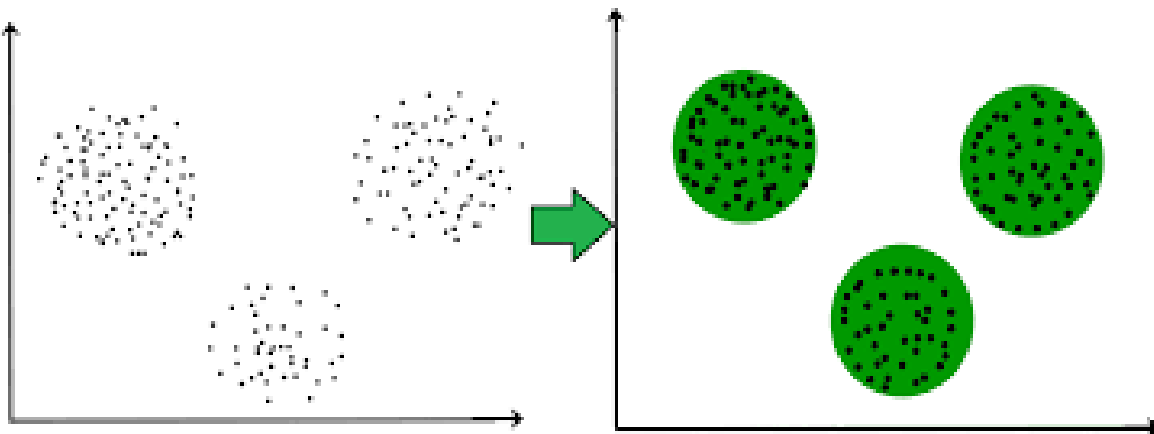
## 1. Introduction

In the field of data science and machine learning, **clustering** and **dimensionality reduction** are two fundamental techniques used for data analysis and preprocessing, particularly in the context of unsupervised learning. These methods help in uncovering hidden patterns in data and in reducing the complexity of high-dimensional datasets, making them more manageable and easier to visualize.

## 2. Body

### Clustering

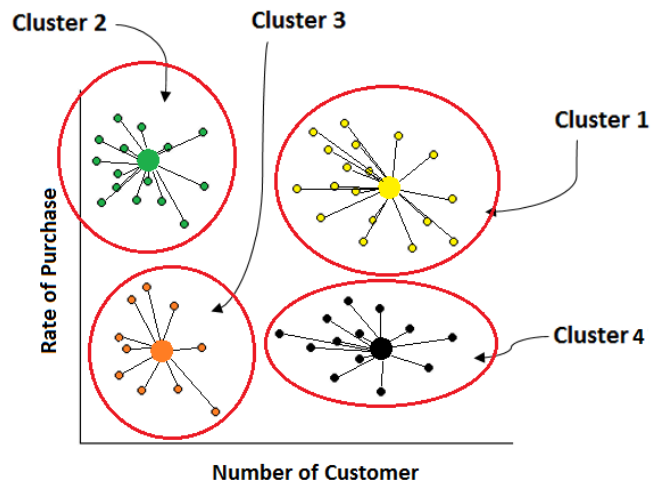
Clustering is a technique used to group a set of objects or data points in such a way that those within the same group (called a cluster) are more similar to each other than to those in other groups. The similarity between data points is usually determined based on a predefined distance metric such as Euclidean distance, Manhattan distance, or cosine similarity. Clustering is often used in various applications, such as customer segmentation, image analysis, pattern recognition, and anomaly detection.



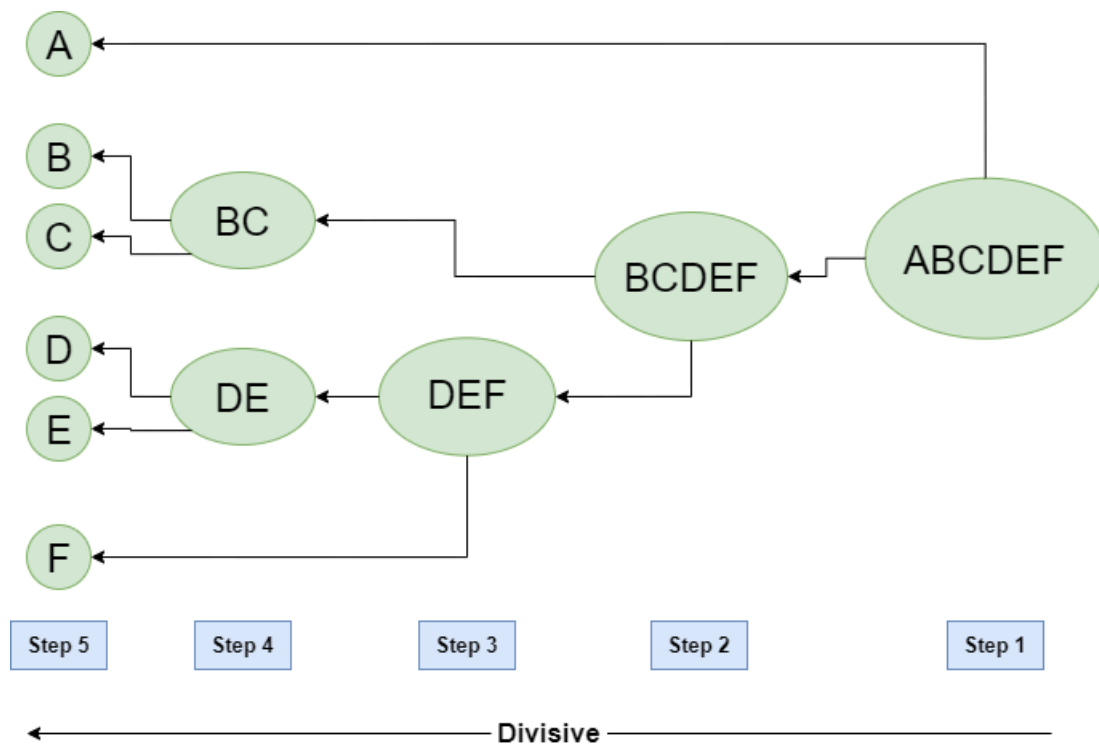
There are several clustering algorithms, each with its own advantages and limitations. Some of the most commonly used clustering methods include:

1. **K-Means Clustering:** This is one of the simplest and most popular clustering algorithms. It partitions the data into  $k$  clusters, where each data point belongs to the cluster with the

nearest mean. The algorithm iteratively refines the clusters by minimizing the variance within each cluster.



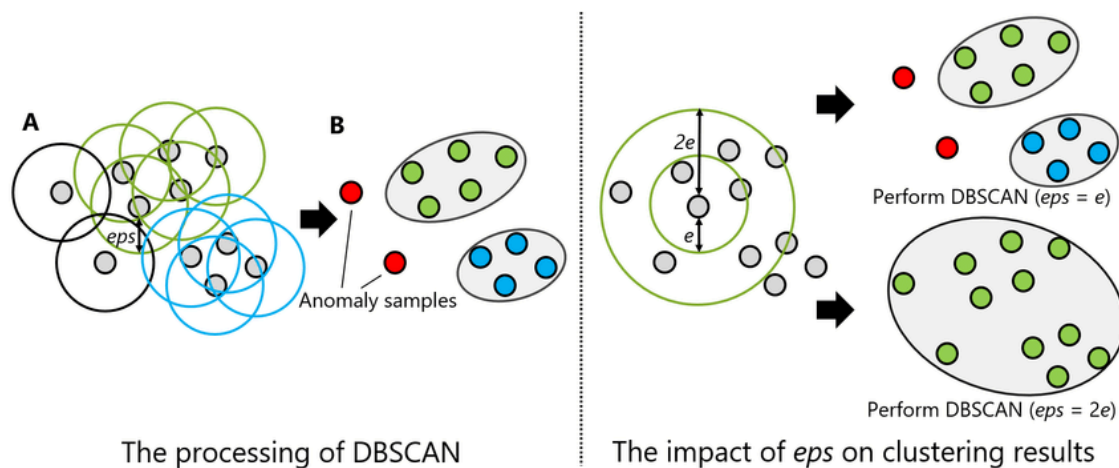
2. **Hierarchical Clustering:** This method builds a hierarchy of clusters either by progressively merging smaller clusters into larger ones (agglomerative) or by splitting larger clusters into smaller ones (divisive). The result is usually represented as a dendrogram, which shows the



relationships between clusters.

3. **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** Unlike K-Means, DBSCAN does not require the number of clusters to be specified in advance. It groups together points that are closely packed together (points with many nearby neighbors) and marks points that lie alone in low-density regions (outliers).

Each of these algorithms has its use cases and is chosen based on the nature of the data and the specific problem at hand.

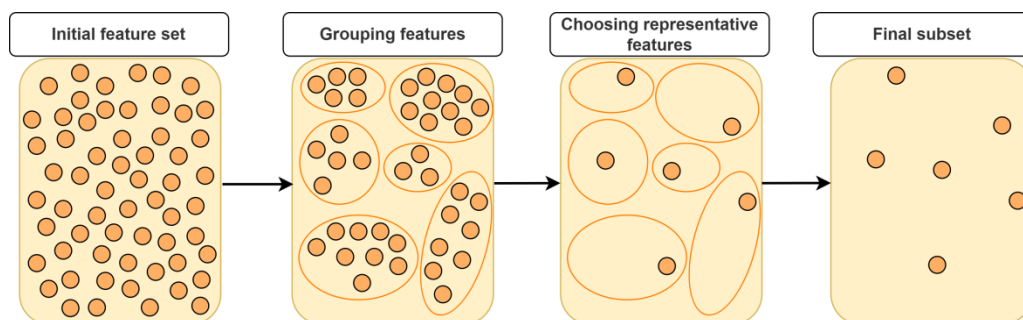


## Dimensionality Reduction

Dimensionality reduction is a technique used to reduce the number of input variables or features in a dataset while retaining as much of the original information as possible. High-dimensional data can be challenging to work with due to the "curse of dimensionality," which can lead to overfitting, increased computational cost, and difficulty in visualizing the data. Dimensionality reduction techniques help to simplify models, reduce noise, and improve interpretability.

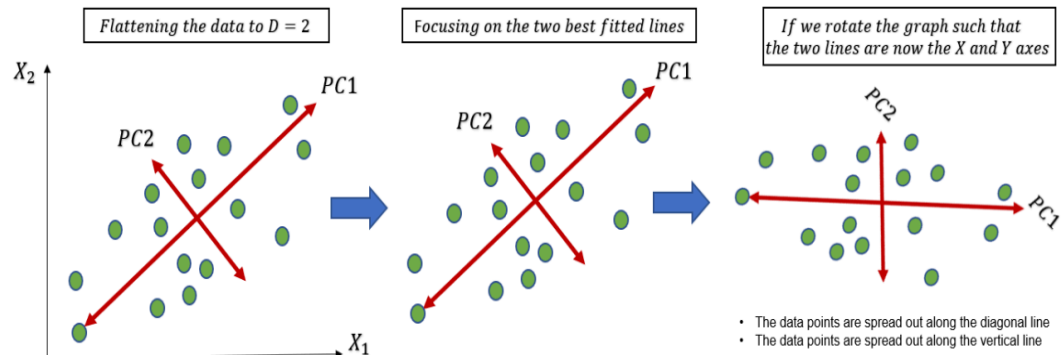
There are two main types of dimensionality reduction techniques:

1. **Feature Selection:** This approach involves selecting a subset of the most important features based on statistical methods, such as correlation coefficients or feature importance scores derived from models like decision trees or random forests.

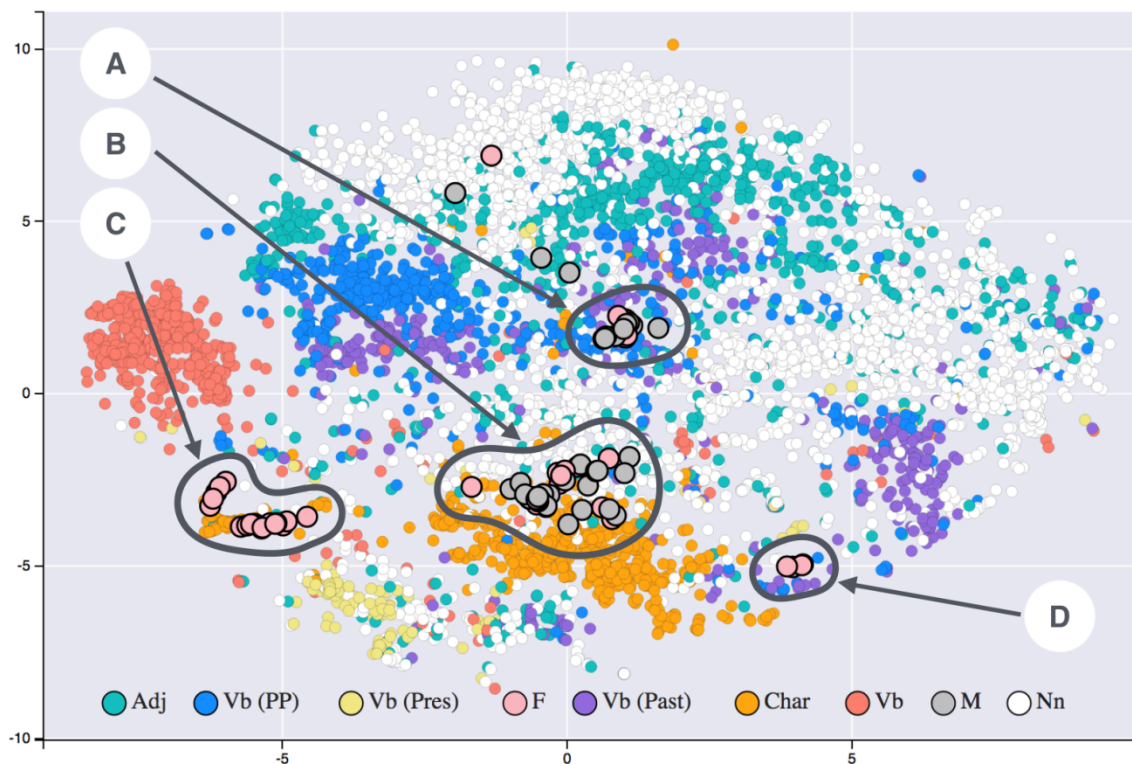


2. **Feature Extraction:** This approach involves transforming the data from a high-dimensional space to a lower-dimensional one. Popular methods for feature extraction include:

- **Principal Component Analysis (PCA):** PCA is a linear dimensionality reduction technique that transforms the data into a new coordinate system where the greatest variance comes to lie on the first few coordinates, known as principal components



**t-Distributed Stochastic Neighbor Embedding (t-SNE):** t-SNE is a non-linear dimensionality reduction technique particularly well-suited for visualizing high-dimensional datasets. It reduces the dimensions by minimizing the divergence between two probability distributions, one representing pairwise similarities in the original space and the other in the reduced space.



### 3. Conclusion

Clustering and dimensionality reduction are crucial techniques in the field of data science, especially for tasks involving unsupervised learning. Clustering helps to identify patterns and natural groupings in data, enabling better understanding and decision-making in various domains such as marketing, biology, and image analysis. On the other hand, dimensionality reduction simplifies complex datasets by reducing the number of variables, which aids in visualization, noise reduction, and improved model performance. Both techniques complement each other, as clustering can benefit from reduced dimensions to operate more efficiently, and dimensionality reduction can leverage clustering to enhance feature extraction. Mastering these techniques allows data scientists to handle complex datasets more effectively and derive meaningful insights from them.

### 4. References

- 📖 Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- 📖 Han, J., Kamber, M., & Pei, J. (2012). *Data Mining: Concepts and Techniques*. Elsevier.
- 📖 Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- 📖 Jolliffe, I. T., & Cadima, J. (2016). "Principal component analysis: a review and recent developments." *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2065).
- 📖 van der Maaten, L., & Hinton, G. (2008). "Visualizing Data using t-SNE." *Journal of Machine Learning Research*, 9(Nov), 2579-2605.