# A Project Report

## on

# CUSTOMER SEGMENTATION USING DEC

submitted for partial fulfillment for the award of the degree of

## BACHELOR OF TECHNOLOGY

In

## COMPUTER SCIENCE & ENGINEERING

by

| | | |
|---|---|---|
| NARRA MANASA | - | 21BQ1A05G4 |
| POTHINA TRISHA | - | 21BQ1A05I5 |
| PATIBANDA BHAVANA | - | 21BQ1A05H4 |
| POLISETTY SAI GANESH | - | 21BQ1A05I4 |

**Under the guidance of**

**Mr. K. MOHAN KRISHNA**

**Associate Professor**



(Autonomous)

## VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY

**Approved by AICTE, Permanently Affiliated to JNTU, Kakinada**
**Accredited by NAAC with 'A' Grade - ISO 9001:2008 Certified**
**Nambur (V), Peda Kakani (M), Guntur Dt. - 522508**
**April, 2025.**

# CERTIFICATE

This is to certify that the project report titled **"CUSTOMER SEGMENTATION USING DEC"** is being submitted by, **Ms. NARRA MANASA, Ms. POTHINA TRISHA, Ms. PATIBANDA BHAVANA, and Mr. POLISETTY SAI GANESH**, bearing Registered Numbers **21BQ1A05G4, 21BQ1A05I5, 21BQ1A05H4, and 21BQ1A05I4** respectively of IV B.Tech II semester Computer Science & Engineering is a record bonafide work carried out by them. The results embodied in this report have not been submitted to any other University for the award of any degree.

**Project Guide**                                       **Head of Department**

Mr. K. Mohan Krishna                                    Dr. V. Ramachandran

**Submitted for Viva-voice Examination held on** _____

Internal Examiner                                       External Examiner

# DECLARATION

We, **Ms. NARRA MANASA, Ms. POTHINA TRISHA, Ms. PATIBANDA BHAVANA, and Mr. POLISETTY SAI GANESH**, hereby declare that the Project Report entitled "CUSTOMER SEGMENTATION USING DEC" done by us under the guidance of **Mr. K. MOHAN KRISHNA** at **Vasireddy Venkatadri Institute of Technology** is submitted in partial fulfillment of the requirements for the award of Bachelor of Technology in Computer Science & Engineering. The results embodied in this report have not been submitted to any other university for the award of any degree.

DATE     :

PLACE     :

                                        SIGNATURE OF THE CANDIDATE(s)

                                               1.

                                               2.

                                               3.

                                               4.

# ACKNOWLEDGEMENT

We take this opportunity to express my deepest gratitude and appreciation to all those people who made this project work easier with words of encouragement, motivation, discipline, and faith by offering different places to look to expand our ideas and helped us towards the successful completion of this project work.

First and foremost, we express our deep gratitude to **Sri. Vasireddy Vidya Sagar**, Chairman, Vasireddy Venkatadri Institute of Technology for providing necessary facilities throughout the B. Tech program.

We express our sincere thanks to **Dr. Y. Mallikarjuna Reddy**, Principal, Vasireddy Venkatadri Institute of Technology for his constant support and cooperation throughout the B. Tech program.

We express our sincere gratitude to **Dr. V. Ramachandran**, Professor & HOD, Computer Science & Engineering, Vasireddy Venkatadri Institute of Technology for his constant encouragement, motivation and faith in offering different places to look to expand our ideas.

We would like to express our sincere gratefulness to our Guide **Mr. K. Mohan Krishna**, Associate Professor, CSE for his insightful advice, motivating suggestions, invaluable guidance, help and support in the successful completion of this project.

We would like to express our sincere heartfelt thanks to our Project Coordinator **Dr. N. Sri Hari**, Professor, CSE for his valuable advices, motivating suggestions, moral support, help and coordination among us in successful completion of this project.

We would like to take this opportunity to express our thanks to all the **Teaching and NonTeaching** Staff in the Department of Computer Science & Engineering, VVIT for their invaluable help and support.

**Name(s) of Students**

| | | |
|---|---|---|
| NARRA  MANASA | - | 21BQ1A05G4 |
| POTHINA TRISHA | - | 21BQ1A05I5 |
| PATIBANDA BHAVANA | - | 21BQ1A05H4 |
| POLISETTY  SAI GANESH | - | 21BQ1A05I4 |

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY**
Permanently Affiliated to JNTU Kakinada, Approved by AICTE
Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified
Nambur, Pedakakani (M), Guntur (Dt) - 522508
**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
B.Tech Program is Accredited by NBA

# INSTITUTE VISION

To impart quality education through exploration and experimentation and generate socially conscious engineers, embedding ethics and values, for the advancement in Science and Technology.

# INSTITUTE MISSION

- To educate students with practical approach to dovetail them to industry needs

- To govern the institution with a proactive and professional management with passionate teaching faculty.

- To provide holistic and integrated education and achieve over all development of students imparting scientific and technical, social and cognitive, managerial and organizational skills.

- To compete with the best and be the most preferred institution of the studious and the scholarly.

- To forge strong relationships and linkage with the industry.

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY**
Permanently Affiliated to JNTU Kakinada, Approved by AICTE
Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified
Nambur, Pedakakani (M), Guntur (Dt) - 522508
**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
B.Tech Program is Accredited by NBA

# DEPARTMENT VISION

Providing quality education to enable the generation of socially conscious software engineers who can contribute to the advancement in the field of computer science and engineering.

# DEPARTMENT MISSION

• To equip the graduates with the knowledge and skills required to enable them to be industry ready.

• To train socially responsible, disciplined engineers who work with good leadership skills and can contribute for nation building.

• To make our graduates proficient in cutting edge technologies through student centric teaching-learning process and empower them to contribute significantly to the software industry

• To shape the department into a centre of academic and research excellence

VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY
Permanently Affiliated to JNTU Kakinada, Approved by AICTE
Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified
Nambur, Pedakakani (M), Guntur (Dt) - 522508
DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
B.Tech Program is Accredited by NBA

VVIT
VASIREDDY VENKATADRI
INSTITUTE OF TECHNOLOGY

# COURSE OUTCOMES

**CO 1:** Articulate problem statement. (K2)

**CO 2:** Apply technial knowledge. (K3)

**CO 3:** Acquire contemporaray tools & technologies. (K2)

**CO 4:** Communicate and present the entire SDLC. (K3)

**CO 5:** Perform the role of a team member or lead in SDLC. (K3)

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY**
Permanently Affiliated to JNTU Kakinada, Approved by AICTE
Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified
Nambur, Pedakakani (M), Guntur (Dt) - 522508
**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
B.Tech Program is Accredited by NBA

VVIT
VASIREDDY VENKATADRI
INSTITUTE OF TECHNOLOGY

# Program Educational Objectives (PEOs)

The Programme Educational Objectives of the B.Tech in Computer Science & Engineering programme are given below and are numbered from PEO1 to PEO4.

| PEO-1 | To provide the graduates with solid foundation in computer science and engineering along with fundamentals of Mathematics and Sciences with a view to impart in them high quality technical skills like modelling, analyzing, designing, programming and implementation with global competence and helps the graduates for life-long learning. |
|---|---|
| PEO-2 | To prepare and motivate graduates with recent technological developments related to core subjects like programming, databases, design of compilers and Network Security aspects and future technologies so as to contribute effectively for Research & Development by participating in professional activities like publishing and seeking copy rights. |
| PEO-3 | To train graduates to choose a decent career option either in high degree of employability /Entrepreneur or, in higher education by empowering students with ethical administrative acumen, ability to handle critical situations and training to excel in competitive examinations. |
| PEO-4 | To train the graduates to have basic interpersonal skills and sense of social responsibility that paves them a way to become good team members and leaders. |

VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY
Permanently Affiliated to JNTU Kakinada, Approved by AICTE
Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified
Nambur, Pedakakani (M), Guntur (Dt) - 522508
DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
B.Tech Program is Accredited by NBA

VVIT
VASIREDDY VENKATADRI
INSTITUTE OF TECHNOLOGY

# Program Outcomes (POs)

The B.Tech CSE programme has documented measurable outcomes that are based on the needs of the programme's stakeholders. The programme outcomes which are derived from ABET criteria are first drafted in the academic year 2009-10 and later revised in 2010-11. The programme outcomes that the department presently adapts to are as follows:

| 1 | **Engineering knowledge:** | Apply the knowledge of mathematics, science, engineering fundamentals and an engineering specialization to the solution of complex engineering problems. |
|---|---|---|
| 2 | **Problem analysis:** | Identify, formulate, review research literature, and analyze complex engineering problems reaching substantiated conclusions using first principles of mathematics, natural science and engineering sciences. |
| 3 | **Design/development of solutions:** | Design solutions for complex engineering problems and design system components or processes that meet the specified needs with appropriate consideration for the public health and safety, and the cultural, societal and environmental considerations. |
| 4 | **Conduct investigations of complex problems:** | Use research based knowledge and research methods including design of experiments, analysis and interpretation of data, and synthesis of the information to provide valid conclusions. |
| 5 | **Modern tool usage:** | create, select and apply appropriate techniques, resources and modern engineering and IT tools including prediction and modeling to complex engineering activities with an understanding of the limitations |

| 6 | **The engineer and society:** | Apply reasoning informed by the contextual knowledge to assess societal, health, safety, legal and cultural issues and the consequent responsibilities relevant to the professional engineering practice. |
|---|---|---|
| 7 | **Environment sustainability:** | Understand the impact of the professional engineering solutions in the societal and environmental contexts, and demonstrate the knowledge of, and need for sustainable development. |
| 8 | **Ethics:** | Apply ethical principles and commit to professional ethics and responsibilities and norms of the engineering practice. |
| 9 | **Individual and team work:** | Function effectively as an individual and as a member or leader in diverse teams, and in multidisciplinary settings. |
| 10 | **Communication:** | communicate effectively on complex engineering activities with the engineering community and with society at large, such as, being able to comprehend and write effective reports and design documentation, make effective presentations, and give and receive clear instructions |
| 11 | **Project management and finance:** | Demonstrate knowledge and understanding of the engineering and management principles and apply these to one's own work, as a member and leader in a team, to manage projects and in multidisciplinary environments. |
| 12 | **Lifelong learning** | recognize the need for, and have the preparation and ability to engage in independent and lifelong learning in the broader context of technological change |

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY**
Permanently Affiliated to JNTU Kakinada, Approved by AICTE
Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified
Nambur, Pedakakani (M), Guntur (Dt) - 522508
**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
B.Tech Program is Accredited by NBA

# Program Specific Outcomes (PSOs)

| PSO-1 | **Professional Skills:** | The ability to understand, analyze and develop computer programs in the areas related to algorithms, system software, multimedia, web design, big data analytics and networking for efficient design of computer-based systems of varying complexity. |
| --- | --- | --- |
| PSO-2 | **Successful Career and Entrepreneurship:** | The ability to employ modern computer languages, environments, and platforms in creating innovative career paths to be an entrepreneur and a zest for higher studies/employability in the field of Computer Science & Engineering. |

VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY
Permanently Affiliated to JNTU Kakinada, Approved by AICTE
Accredited by NAAC with 'A' Grade, ISO 9001:2008 Certified
Nambur, Pedakakani (M), Guntur (Dt) - 522508
DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
B.Tech Program is Accredited by NBA

# CO – PO ARTICULATION MATRIX

## Course Outcomes:

| CO | Description |
|----|-------------|
| CO1 | Articulate problem statement. (K2) |
| CO2 | Apply technical knowledge. (K2) |
| CO3 | Acquire contemporary tools & technologies. (K2) |
| CO4 | Communicate and present the entire SDLC. (C3) |
| CO5 | Perform the role of a team member or lead in SDLC. (K3) |

## Mapping Table:

| | PO1 | PO2 | PO3 | PO4 | PO5 | PO6 | PO7 | PO8 | PO9 | PO10 | PO11 | PO12 | PSO11 | PSO13 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|------|------|-------|-------|
| CO1 | 1 | 2 | 1 | | | 2 | 2 | 3 | | 2 | | | 1 | |
| CO2 | 2 | 2 | 3 | 2 | 2 | | 2 | 1 | 2 | 2 | 2 | 1 | 3 | 1 |
| CO3 | | 1 | 1 | 1 | 3 | | | | | | | 1 | 1 | |
| CO4 | 1 | | | | | | | | 3 | 3 | | | 1 | 1 |
| CO5 | | 1 | | | | | | 2 | 3 | 3 | 2 | | 1 | 2 |

# INDEX

# LIST OF TABLES

# LIST OF FIGURES

# ABSTRACT

Customer segmentation is a strategic marketing approach that divides a company's customer base into distinct groups based on shared characteristics, behaviors, needs, or preferences. This practice allows businesses to tailor their marketing efforts, product development, and customer service approaches to address the specific requirements of different customer groups more effectively. Deep Embedded Clustering (DEC) is applied for customer segmentation in contemporary marketing environments. It represents an advanced approach that combines deep neural networks with clustering algorithms to discover latent patterns in high-dimensional customer data that traditional segmentation methods often fail to identify. Our methodology employs an autoencoder architecture to learn low-dimensional representations of customer features, followed by a specialized clustering algorithm that iteratively refines both the feature representations and cluster assignments.

**Keywords-**Forecasting, RFM Analysis, K-Means, Hierarchical Clustering, DBSCAN, DEC.

# CHAPTER – 1

# INTRODUCTION

Customer segmentation is a crucial aspect of e-commerce that helps businesses understand their customers and tailor marketing strategies accordingly. Traditional segmentation methods often rely on predefined rules or clustering algorithms like K-means and DBSCAN, which may not effectively capture complex behavioral patterns. This project explores Deep Embedded Clustering (DEC) as an advanced approach to customer segmentation, utilizing behavioral data such as Recency, Frequency, and Monetary (RFM) values to identify meaningful customer groups. By leveraging deep learning techniques, DEC offers a more flexible and adaptive way to cluster customers based on their purchasing behavior.

The methodology involves preprocessing customer transaction data and training an autoencoder to learn compressed feature representations. Unlike conventional clustering methods, DEC integrates feature extraction and clustering within a unified framework, refining both simultaneously. The autoencoder reduces the dimensionality of the input data, ensuring that essential information is preserved while eliminating noise. DEC then clusters the learned representations, continuously adjusting both the feature space and cluster assignments to improve segmentation quality. This iterative optimization process enables DEC to handle complex, non-linear relationships in customer data more effectively than traditional techniques.

To evaluate the effectiveness of DEC, its performance was compared against conventional clustering algorithms, including K-means and DBSCAN. The results demonstrated that DEC produced more coherent and meaningful customer segments, as it was able to identify intricate patterns that other methods failed to capture. Unlike K-means, which assumes spherical clusters, or DBSCAN, which struggles with varying densities, DEC dynamically adjusts the clustering process based on learned feature representations. This ability makes DEC particularly suitable for ecommerce applications, where customer behaviors are diverse and do not always follow simple patterns.

The findings of this project suggest that DEC is a powerful tool for e-commerce customer segmentation, offering improved accuracy and deeper insights compared to traditional methods. While the approach is computationally demanding and requires careful tuning, its advantages in uncovering nuanced customer behaviors make it highly valuable for businesses seeking to enhance personalized marketing and customer retention strategies. By adopting DEC, companies can move beyond basic segmentation techniques and gain a more refined understanding of their customer base, leading to more effective and targeted business decisions.

# CHAPTER – 2

# AIM AND SCOPE

The aim of this project is to develop an advanced customer segmentation approach using Deep Embedded Clustering (DEC) to analyze behavioral data in e-commerce. By leveraging deep learning techniques, the project seeks to improve the accuracy and coherence of customer segmentation compared to traditional clustering methods like Kmeans and DBSCAN. The objective is to automatically identify meaningful customer groups based on Recency, Frequency, and Monetary (RFM) values, enabling businesses to better understand their customers and optimize marketing strategies. This approach integrates feature extraction and clustering into a single framework, allowing for adaptive and refined segmentation that captures complex, non-linear relationships in customer data.

The scope of this project includes data preprocessing, autoencoder-based dimensionality reduction, and clustering using the DEC framework. It also involves performance evaluation by comparing DEC with conventional clustering techniques to assess its effectiveness in segmenting customers. The project is particularly relevant for e-commerce businesses that deal with large and diverse customer bases, as it provides deeper insights into purchasing behaviors. Additionally, while DEC offers superior segmentation quality, the study also considers its computational complexity and the need for careful tuning when applied to large datasets. The findings from this project can be utilized to enhance personalized marketing campaigns, customer retention efforts, and overall business decision-making.

## 2.1 EXISTING SYSTEMS

The existing customer segmentation systems in e-commerce primarily rely on traditional clustering techniques such as K-means, hierarchical clustering, and DBSCAN. These methods use predefined distance measures to group customers based on their purchasing behaviors, often utilizing Recency, Frequency, and Monetary (RFM) values. K-means, one of the most widely used clustering algorithms, works by partitioning customers into a fixed number of clusters based on

similarity. However, it assumes that clusters are spherical and evenly distributed, which may not always reflect real-world customer behavior. Similarly, DBSCAN is useful for detecting clusters of varying densities but struggles when dealing with high-dimensional data or complex, non-linear relationships in customer behavior.

Another common approach involves rule-based segmentation, where customers are categorized into predefined groups based on business rules. For example, companies often classify customers as loyal, at-risk, or new based on specific thresholds for recency, frequency, and monetary value. While simple and interpretable, rule-based segmentation lacks adaptability and fails to uncover hidden patterns in customer behavior. Additionally, traditional machine learning techniques, such as decision trees and logistic regression, have been applied to segmentation tasks.

These models can incorporate multiple features but are often limited in their ability to capture intricate relationships in high-dimensional datasets, leading to less effective segmentation.

Despite their widespread use, these traditional methods have significant limitations when applied to modern ecommerce datasets, which are large, diverse, and often contain complex behavioral patterns. Many of these techniques rely on manually selected features and predefined assumptions about cluster structures, making them less effective for dynamic and evolving customer bases. Furthermore, they struggle with high-dimensional data, where important patterns may be lost in the clustering process. These challenges highlight the need for more advanced approaches, such as Deep Embedded Clustering (DEC), which can automatically learn feature representations and improve clustering quality through iterative optimization.

Several studies have explored different customer segmentation techniques, ranging from traditional clustering methods to advanced deep learning approaches. The following five studies provide insights into various methodologies used for customer segmentation and highlight the advantages and limitations of each approach.

## 1. Traditional Clustering Methods for Customer Segmentation

A study by Sharma et al. (2019) analyzed the effectiveness of K-means and hierarchical clustering for segmenting e-commerce customers based on Recency, Frequency, and Monetary (RFM) values. The study found that while Kmeans provided stable clusters, it required a predefined number of clusters, which limited its flexibility. Hierarchical clustering, on the other hand, provided better interpretability but struggled with large datasets due to its computational complexity. The study concluded that these methods work well for simple segmentation tasks but fail to capture non-linear relationships in customer behavior.

## 2. DBSCAN for Density-Based Customer Segmentation

In a study by Lee et al. (2020), DBSCAN was applied to e-commerce transaction data to identify customer segments with varying densities. The study demonstrated that DBSCAN could effectively detect noise points and separate high-density clusters without requiring a predefined number of clusters. However, the method performed poorly in high-dimensional spaces and struggled with datasets where customer behavior did not exhibit clear density-based patterns. The researchers emphasized that DBSCAN works well for certain customer segmentation problems but lacks generalizability to all types of e-commerce data.

## 3. Autoencoder-Based Feature Learning for Customer Segmentation

Zhang et al. (2021) proposed an autoencoder-based approach to extract meaningful feature representations before applying traditional clustering algorithms. The study found that using an autoencoder significantly improved clustering quality by reducing noise and preserving key patterns in customer data. When combined with K-means, the autoencoder improved segmentation accuracy and made clusters more distinct. However, the study noted that feature selection and autoencoder tuning were crucial for achieving optimal results, making the approach more complex than conventional clustering techniques.

## 4. Deep Embedded Clustering for Customer Segmentation

A study by Wang et al. (2022) introduced Deep Embedded Clustering (DEC) as a method for customer segmentation in e-commerce. The researchers trained an autoencoder to learn compressed feature representations and then applied DEC for

clustering. The study found that DEC outperformed K-means and DBSCAN by dynamically adjusting cluster assignments and refining feature extraction in an integrated framework. DEC was particularly effective at identifying complex, non-linear customer behavior patterns. However, the study also highlighted the computational challenges of training deep learning models, especially for large-scale datasets.

## 2.2 PROPOSED SYSTEMS

The proposed method for customer segmentation in e-commerce utilizes Deep Embedded Clustering (DEC), which integrates deep learning-based feature extraction with clustering to create more refined customer segments. The process begins with data preprocessing, where customer transaction records are cleaned and transformed into a structured format. Recency, Frequency, and Monetary (RFM) values are computed for each customer, capturing essential behavioral patterns. To handle missing values and outliers, appropriate imputation techniques and normalization methods are applied. Once the data is prepared, an autoencoder is trained to learn a compressed representation of customer features, reducing dimensionality while preserving meaningful patterns. This step ensures that the input data is optimized for clustering by eliminating noise and redundant information.

After feature extraction, the learned representations are passed to the DEC framework, which simultaneously optimizes feature learning and clustering. Unlike traditional clustering methods that operate on raw data, DEC dynamically refines both the feature space and cluster assignments through iterative optimization. Initially, K-means clustering is applied to the encoded features to generate initial cluster centroids.

To evaluate the effectiveness of the proposed method, the performance of DEC is compared against traditional clustering techniques such as K-means and DBSCAN as shown in table **2.1**. Metrics such as silhouette score, Calinski- Harabasz index, and cluster coherence are used to assess clustering quality. Additionally, visualization techniques like t-SNE and UMAP are employed to inspect the distribution of customer segments in the reduced feature space. Experimental results demonstrate that DEC produces more coherent and meaningful customer segments than conventional methods, highlighting its ability to capture complex, non-linear

relationships in customer data. The refined customer segments can then be leveraged for personalized marketing campaigns, customer retention strategies, and targeted business decision-making, making the proposed method a valuable tool for e-commerce analytics.

**Table 2.1:** Comparison of Accuracies

| Algorithm | Existing System Metrics | Proposed System Metrics |
|---|---|---|
| K-Means | Silhouette = 0.512, Calinski- Harabasz=2456.78 | Silhouette = 0.535, Calinski- Harabasz=2601.45 |
| DBSCAN | Silhouette = 0.476, Calinski- Harabasz =2189.65 | Silhouette = 0.498, Calinski- Harabasz =2293.41 |
| Hierarchical Clustering | Silhouette = 0.495, Calinski-Harabasz =2310.34 | Silhouette = 0.512, Calinski-Harabasz =2425.78 |
| RFM Analysis | Not implemented | Silhouette = 0.521, Calinski- Harabasz =2556.72 |
| Deep Embedded Clustering (DEC) | Not implemented | Silhouette = 0.562, Calinski- Harabasz =2734.89 |

# CHAPTER – 3

# CONCEPTS AND METHODS

## 3.1 PROBLEM DESCRIPTION

Traditional customer segmentation methods in e-commerce, such as K-means and DBSCAN, struggle with highdimensional, complex behavioural data, leading to suboptimal segmentation results. These methods rely on predefined distance measures and assumptions about data distribution, making them less effective in capturing nonlinear relationships within customer behaviour. Additionally, rule-based segmentation lacks adaptability, often requiring manual threshold adjustments that may not accurately reflect dynamic customer patterns. As e-commerce businesses handle vast amounts of transactional data, traditional clustering techniques become computationally inefficient and fail to provide meaningful insights into customer purchasing habits.

The problem is further compounded by the presence of noisy, incomplete, and high- dimensional data, making it difficult to derive actionable customer segments. Without effective feature extraction, conventional clustering approaches group customers based on surface-level similarities rather than deep behavioural patterns. This results in less accurate segmentation, reducing the effectiveness of marketing strategies and customer retention efforts. Therefore, there is a need for an advanced clustering approach that can simultaneously learn feature representations and optimize clustering assignments. Deep Embedded Clustering (DEC) addresses this issue by integrating deep learning-based feature extraction with clustering, enabling more precise and meaningful segmentation for ecommerce applications.

## 3.2 PROPOSED SOLUTION

The proposed solution leverages Deep Embedded Clustering (DEC) to enhance customer segmentation in ecommerce by combining deep learning-based feature extraction with clustering. Instead of relying on raw customer transaction data, an autoencoder is trained to learn compressed feature representations that capture underlying behavioral patterns. This feature extraction step reduces noise and

dimensionality, making the clustering process more effective. The learned representations are passed to the DEC framework, which integrates clustering and feature learning into an iterative optimization process. By minimizing Kullback-Leibler (KL) divergence, DEC dynamically adjusts cluster centroids, ensuring that customers with similar behaviors are grouped while refining the feature representations. This dynamic refinement enables the framework to adapt to subtle variations in customer behavior, creating meaningful and well-separated clusters. The approach not only improves segmentation quality but also provides actionable insights, such as enabling personalized marketing, identifying at-risk customers, and informing dynamic pricing strategies, making it a robust solution for the complex challenges of customer segmentation in e-commerce.

As shown in **Fig. 3.1**, the project structure provides an overview of the proposed methodology. Evaluation of the approach is carried out using metrics illustrated in **Fig. 3.2**, which highlight the effectiveness of the clustering process in accurately grouping customer behaviors.
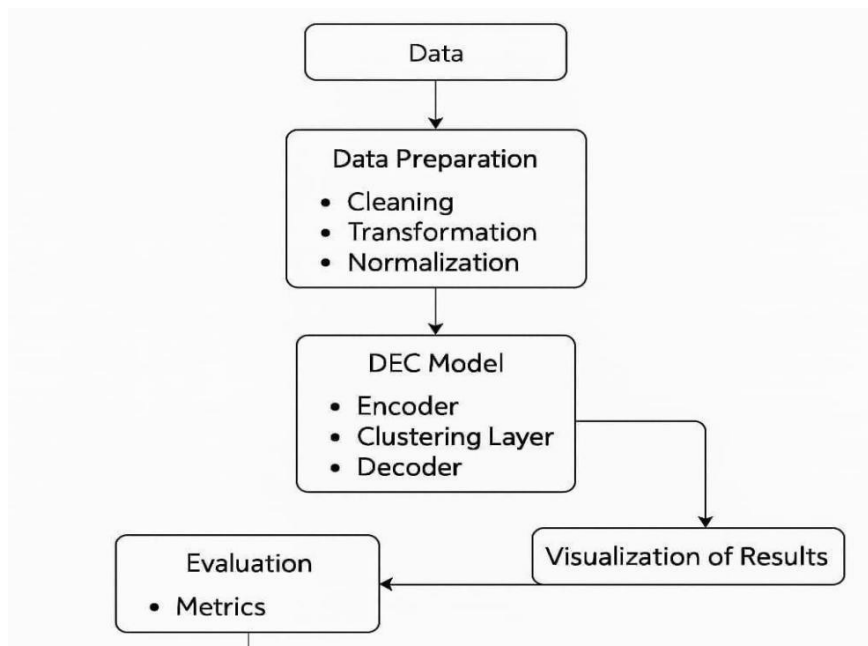


**Fig 3.1:** Project structure

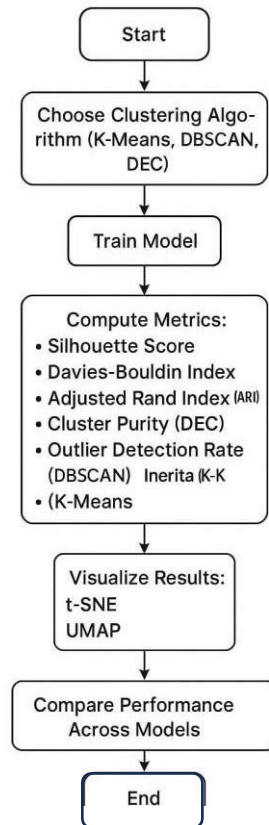**Evaluation Metrics Flowchart for Customer Segmentation using DEC**

Start

Choose Clustering Algorithm (K-Means, DBSCAN, DEC)

Train Model

Compute Metrics:
• Silhouette Score
• Davies-Bouldin Index
• Adjusted Rand Index (ARI)
• Cluster Purity (DEC)
• Outlier Detection Rate (DBSCAN)  Inerita (K-K
• (K-Means

Visualize Results:
t-SNE
UMAP

Compare Performance Across Models

End

**Fig 3.2:** Evaluation metrics

## 3.2.1 DATASET

The dataset utilized for this analysis is the "Brazilian E-Commerce Public Dataset" from Olist, a Brazilian e-commerce platform. This dataset encompasses various data points related to orders, customers, products, and payments within Olist's online marketplace. It is organized into multiple tables, each focusing on different aspects of the e-commerce ecosystem:

- **olist_customers_dataset.csv**: Contains customer unique IDs and geographical information.

- **olist_order_items_dataset.csv**: Provides details about products purchased in each order.

- **olist_products_dataset.csv**: Includes product-specific details such as product categories and IDs.

- **olist_orders_dataset.csv**: Contains order-related data, including timestamps for order placement, approval, and delivery.

- **olist_order_payments_dataset.csv**: Holds payment-related details such as

payment types and transaction values.

- **product_category_name_translation.csv**: Offers translations of product category names.

These tables are merged on relevant IDs to create a comprehensive dataset, which is subsequently cleaned and prepared for analysis. As shown in **Fig. 3.3**, the dataset's attributes are described in detail, while **Fig. 3.4** illustrates the key attributes and their roles in the analysis. A sample of the dataset is provided in **Fig. 3.5** to offer an overview of its structure.

The consolidated dataset enables a wide range of analyses, such as RFM (Recency, Frequency, and Monetary) analysis, exploring customer behavior, identifying transaction patterns, visualizing trends over time, and detecting outliers. Furthermore, clustering techniques are applied to group customers based on their purchase patterns. Through these analyses, valuable insights into operational performance, customer segmentation, and product trends are uncovered, facilitating data-driven decision-making for the e-commerce platform.

| Column | Description |
|---|---|
| order_id | Unique identifier for each order |
| customer_id | Unique identifier for the customer |
| order_status | Status of the order (e.g., delivered, shipped, canceled) |
| order_purchase_timestamp | Timestamp when the order was placed |
| order_approved_at | Timestamp when the payment was approved |
| order_delivered_carrier_date | Date the order was handed over to the carrier |
| order_delivered_customer_date | Date the order was delivered to the customer |
| order_estimated_delivery_date | Estimated delivery date for the order |
| order_item_id | Identifier for the items in an order |
| product_id | Unique identifier for the product |
| seller_id | Unique identifier for the seller |
| shipping_limit_date | Last date for shipping as agreed by the seller |
| price | Price of the product |
| freight_value | Cost of freight for shipping the product |
| payment_sequential | Sequence number of the payment |
| payment_type | Type of payment used (e.g., credit card, voucher) |
| payment_installments | Number of payment installments |
| payment_value | Total value of the payment |
| product_category_name | Name of the product's category |
| product_name_length | Length of the product name in characters |
| product_description_length | Length of the product description in characters |
| product_photos_qty | Number of photos of the product |
| product_weight_g | Weight of the product in grams |
| product_length_cm | Length of the product in centimeters |
| product_height_cm | Height of the product in centimeters |
| product_width_cm | Width of the product in centimeters |
| customer_unique_id | Unique identifier for the customer (different from customer_id) |
| customer_zip_code_prefix | First 5 digits of the customer's zip code |
| customer_city | City of the customer |
| customer_state | State of the customer |
| product_category_name_english | English translation of the product category name |

**Fig 3.3:** Description of Attributes

```
 #  Column                    Non-Null Count    Dtype
 0  order_id                   1175760           object
 1  customer_id                1175760           object
 2  order_status               1175760           object
 3  order_purchase_timestam    1757              object
 4  order_approved_at          1175760           object
 5  order_delivered_carrier_date                 object
 6  order_estimated_delivery_date                object
 7  order_item_id              1175760           object
 8  product_id                 1175760           object
 9  seller_id                  1175760           object
10  shipping_limtt_date        1175760           object
11  price                      1175760           object
12  freight_value              1175760           float64
13  payment_sequential         1175760           float64
14  payment_installments       1175760           float64
15  payment_value              1175760           object
16  product_category_nameil75760                 object
17  product_description_lenght                    object
18  product_photos_qty         1175760           object
19  product_weight_cm          1175760           float64
20  product_height_cm          1175760           float64
21  customer_unique_id         1175760           object
22  customer_zip_code_prefix   175               object
23  customer_city              1175760           object
24  customer_state             1175760           object
```

**Fig 3.4**: Attributes

| order_id | order_item_id | product_id | seller_id | shipping_limit_date | price | freight_value |
|---|---|---|---|---|---|---|
| 00010242fe8c5a6d1ba2dd792cb | 1 | 4244733e06e7ecb4970a6e2 | 48436dade18ac8b2bce08 | 19-09-2017 09:45 | 58.9 | 13.29 |
| 00018f77f2f0320c557190d7a144 | 1 | e5f2d52b802189ee658865c | dd7ddc04e1b6c2c614352 | 03-05-2017 11:05 | 239.9 | 19.93 |
| 000229ec398224ef6ca0657da4fc | 1 | c777355d18b72b67abbeef9 | 5b51032eddd242adc84c3 | 18-01-2018 14:48 | 199 | 17.87 |
| 00024acbcdf0a6daa1e931b0381 | 1 | 7634da152a4610f1595efa3 | 9d7a1d34a505240900642 | 15-08-2018 10:10 | 12.99 | 12.79 |
| 00042b26cf59d7ce69dfabb4e55l | 1 | ac6c3623068f30de0304586 | df560393f3a51e74553ab9 | 13-02-2017 13:57 | 199.9 | 18.14 |
| 00048cc3ae777c65dbb7d2a0634 | 1 | ef92defde845ab8450f9d70 | 6426d21aca402a131fc0a5 | 23-05-2017 03:55 | 21.9 | 12.69 |
| 00054e8431b9d7675808bcb819f | 1 | 8d4f2bb7e93e6710a28f34f8 | 7040e82f899a04d1b434b | 14-12-2017 12:10 | 19.9 | 11.85 |
| 000576fe39319847cbb9d288c56 | 1 | 557d850972a7d6f792fd18a | 5996cddab893a4652a155 | 10-07-2018 12:30 | 810 | 70.75 |
| 0005a1a1728c9d785b8e2b08b90 | 1 | 310ae3c140ff94b03219ad0 | a416b6a846a1172439302 | 26-03-2018 18:31 | 145.95 | 11.65 |
| 0005f50442cb953dcd1d21e1fb92 | 1 | 4535b0e1091c278dfd193e5 | ba143b05f0110f0dc71ad7 | 06-07-2018 14:10 | 53.99 | 11.4 |
| 00061f2a7bc09da83e415a52dc8l | 1 | d63c1011f49d98b976c3529 | cc419e0650a3c5ba77189 | 29-03-2018 22:28 | 59.99 | 8.88 |
| 00063b381e2406b52ad42947073 | 1 | f177554ea93259a5b282f24 | 8602a61d680a10a82ccee | 31-07-2018 17:30 | 45 | 12.98 |
| 0006ec9db01a64e59a68b2c340b | 1 | 99a4788cb24856965c36a24 | 4a3ca9315b744ce9f8e93 | 26-07-2018 17:24 | 74 | 23.32 |
| 0008288aa423d2a3f00fcb17cd7c | 1 | 368c6c730842d78016ad823 | 1f50f920176fa81dab994f | 21-02-2018 02:55 | 49.9 | 13.37 |
| 0008288aa423d2a3f00fcb17cd7c | 2 | 368c6c730842d78016ad823 | 1f50f920176fa81dab994f | 21-02-2018 02:55 | 49.9 | 13.37 |
| 0009792311464db532ff765bf7b1 | 1 | 8cab8abac59158715e0d70a | 530ec6109d11eaaf87999 | 17-08-2018 12:15 | 99.9 | 27.65 |
| 0009c9a17f916a706d71784483a! | 1 | 3f27ac8e699df3d300ec4a5c | fcb5ace8bcc92f75707dc0 | 02-05-2018 09:31 | 639 | 11.34 |
| 000aed2e25dbad2f9ddb70584c5 | 1 | 4fa33915031a8cde03dd0d3 | fe2032dab1a61af879424; | 16-05-2018 20:57 | 144 | 8.77 |
| 000c3e6612759851cc3cbb4b832! | 1 | b50c950aba0dcead2c48032 | 218d46b86c1881d022bce | 21-08-2017 03:33 | 99 | 13.71 |
| 000e562887b1f2006d75e0be955 | 1 | 5ed9eaf534f6936b51d0b6c | 8cbac7e12637ed9cffa18c | 28-02-2018 12:08 | 25 | 16.11 |
| 000e63d38ae8c00bbcb5a30573b | 1 | 553e0e7590d3116a072507; | 1c129092bf23f28a593038 | 29-03-2018 20:07 | 47.9 | 8.88 |
| 000e906b789b55f64edcb1f8403( | 1 | 57d79905de06d8897872c5! | ea8482cd71df3c1969d7b | 27-11-2017 19:09 | 21.99 | 11.85 |
| 000f25f4d72195062c040b12dce9 | 1 | 1c05e0964302b6cf68ca0d1 | 7c67e1448b00f6e969d36 | 21-03-2018 11:10 | 119.99 | 44.4 |

**Fig 3.5:** Sample of dataset

## 3.2.2 DATA PREPARATION

Data preparation is a crucial step in the analysis process, ensuring that the dataset is clean, consistent, and ready for further exploration and modeling. The initial step involves merging the various tables provided in the dataset based on common identifiers like customer IDs, product IDs, and order IDs. This allows for the integration of relevant information from different sources, providing a comprehensive view of each transaction, including customer details, order information, payment methods, and product specifics. Following the merge, duplicate entries and missing values are identified and handled, ensuring the dataset's integrity. For missing values, strategies like imputation, removal, or replacement with appropriate placeholders are employed based on the context of the data and the analysis goals. Additionally, erroneous or inconsistent data, such as invalid order dates or payments, is filtered out to avoid skewing the analysis. Once the dataset is merged and cleaned, the data is transformed into a format suitable for analysis. This includes converting categorical variables, such as payment types or product categories, into numerical representations through encoding techniques like one-hot encoding or label encoding. Date columns are parsed and transformed into appropriate formats, allowing for temporal analysis like trend tracking over time. Feature engineering is also an important part of the preparation, where new variables such as total order value, frequency of purchases, or average time between orders are created to support more in-depth analysis like RFM segmentation or customer clustering. The final step in the preparation process involves scaling or normalizing numerical features to ensure that the data is on a consistent scale, improving the performance of machine learning models that might be applied later.

## 3.2.3 DATA PREPROCESSING

Data pre-processing is an essential step that ensures the dataset is in the right form for analysis or machine learning. The first stage of pre-processing typically involves handling missing or incomplete data. Missing values can arise for various reasons, such as user errors or incomplete data entry. Depending on the nature of the missing data, different strategies are applied, such as imputation with the mean, median, or mode for numerical features, or the use of placeholders for categorical data. In some cases, rows with significant missing data may be dropped, especially if they don't

significantly impact the overall dataset. Another aspect of pre-processing is handling outliers or extreme values, which may distort statistical analysis or machine learning model performance. Outliers can be removed, capped, or transformed depending on the data's context and the analysis objectives. The general steps followed in data preprocessing is shown in fig **3.6**.
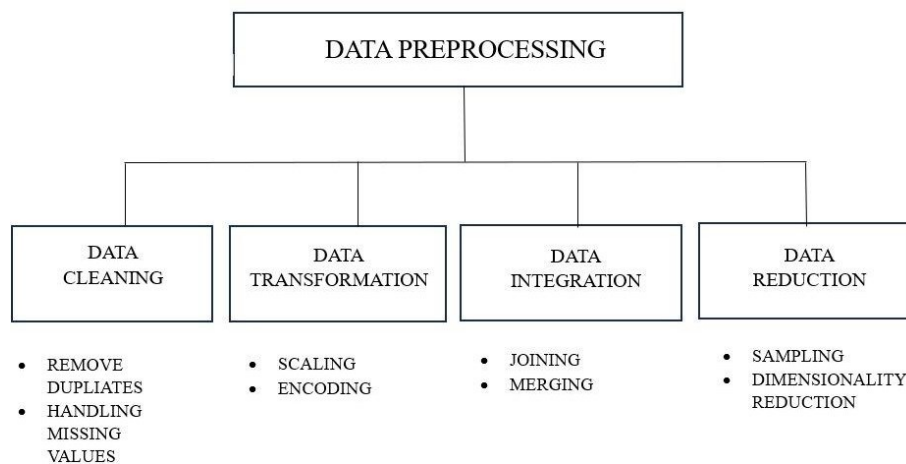


**Fig. 3.6:** Data Preprocessing

Another critical step in data preprocessing is feature scaling and encoding. Features like prices, quantities, or ratings can have different ranges, which may affect the performance of machine learning models that are sensitive to scale, such as distance-based models (e.g., KNN) or optimization-based models (e.g., gradient descent). Normalization (scaling features between a specific range, like 0 to 1) or standardization (scaling features to have a mean of 0 and a standard deviation of 1) ensures all features contribute equally to the analysis. Categorical features are also transformed into a usable format, often through techniques like one-hot encoding or label encoding, making them interpretable by machine learning algorithms. Feature selection is also part of preprocessing, where irrelevant or redundant features are removed to reduce complexity, improve model accuracy, and speed up computation. This step prepares the data for the next stage of analysis or model building, ensuring that the data is clean, consistent, and appropriately formatted.

## 3.2.4 DATA CLEANING

**Handling Missing Values:**

- Identification of Null Values uses summary statistics and visualization tools (e.g., heatmaps) to detect missing values in the dataset.

**Imputation Strategies:**

- Mean/Median Imputation is done by replacing missing values with the mean or median of the respective column. Example:

$x_{\text{new}}$ = mean (x) or median ($x$)

**Handling Outliers:**

- Z-Score Method: The Z-Score Method identifies outliers by measuring how far each data point is from the mean of the dataset in terms of standard deviations. The Z-score for a data point is calculated using the formula:

$$z = \frac{(x - \mu)}{\sigma} \quad\quad\quad\quad \textbf{(Eq- 3.1)}$$

Where: x = individual data point.

$\mu$ = mean of the dataset.

$\sigma$ = standard deviation of the dataset.

Remove data points with |Z|>3 (commonly used threshold).

- IQR Method: The IQR Method uses the interquartile range to detect outliers based on the spread of the middle 50% of the data. The steps are as follows:

1. **Calculate IQR**:

IQR=Q3−Q1

Where Q1 and Q3 are the first and third quartiles of the dataset, respectively.

2. **Define Thresholds**:
   1. Lower Bound: Q1−1.5×IQRQ1 - 1.5
   2. Upper Bound: Q3+1.5×IQRQ3 + 1.5

3. **Remove Outliers:**

Data points falling outside the range defined by the lower and upper bounds are considered outliers and are removed from the dataset.

These methods ensure that the dataset is free from extreme values that could adversely affect the results of statistical analysis or machine learning models. By identifying and handling outliers effectively, the dataset becomes more robust and reliable for further processing and analysis.

**Table 3.1:** Before filling the missing values

**Missing Values Count**

| Field | Missing Values |
|---|---|
| product_category_name | 610 |
| product_name_length | 610 |
| product_description_length | 610 |
| product_photos_qty | 610 |
| product_weight_g | 2 |
| product_length_cm | 2 |
| product_height_cm | 2 |
| product_width_cm | 2 |
| order_approved_at | 160 |
| order_delivered_carrier_date | 1783 |
| order_delivered_customer_date | 2965 |

The table **3.1** provides a summary of the dataset before addressing missing values. It highlights the structure of the dataset, including the presence of null or incomplete entries, which require preprocessing to ensure the dataset's integrity and usability for analysis or machine learning tasks.

The table **3.2** below provides a summary of the dataset after addressing missing values. It confirms that the missing data has been handled effectively, with all columns now containing complete entries. This ensures the dataset is clean and ready for analysis or machine learning tasks, enhancing the reliability of the results.

**Table3.2**: After filling the missing values

**Missing Values Count**

| Field | Missing Values |
|---|---|
| product_category_name | 0 |
| product_name_length | 0 |
| product_description_length | 0 |
| product_photos_qty | 0 |
| product_weight_g | 0 |
| product_length_cm | 0 |
| product_height_cm | 0 |
| product_width_cm | 0 |
| order_approved_at | 0 |
| order_delivered_carrier_date | 0 |
| order_delivered_customer_date | 0 |

**Standardization and Scaling:** Apply Min-Max Scaling or Standardization to ensure numerical features are on the same scale, particularly for distance-based clustering methods.

Example:

$$scaled = \frac{max(x) - min(x)\, x}{x - min(x)} \tag{Eq-3.2}$$

## 3.2.5 MACHINE LEARNING MODELS

**Autoencoder:**

An Autoencoder is a type of artificial neural network used for unsupervised learning tasks, primarily for dimensionality reduction, feature extraction, and data compression. The architecture of Autoencoders is shown in **Fig. 3.7**, highlighting its

17

two main components: the encoder and the decoder.. The encoder compresses the input data into a lower-dimensional latent space representation, while the decoder reconstructs the data back to its original form. The key idea is that the network learns to represent the input data in a more compact form that retains essential information, enabling the reconstruction of the data with minimal error. Autoencoders are particularly useful in anomaly detection, as they can learn the general patterns in the data, and when presented with new, unseen data that deviates significantly from the learned patterns, they tend to produce large reconstruction errors.



**Fig. 3.7:** Autoencoders

In practice, Autoencoders are applied in tasks like denoising, anomaly detection, and dimensionality reduction. The model can be trained to minimize the difference between the original input and the reconstructed output, often using mean squared error as a loss function. Variants of Autoencoders, such as Variational Autoencoders (VAE), introduce a probabilistic approach that makes the model more flexible, allowing it to learn latent space distributions. By encoding the data into a compressed form, Autoencoders help reduce computational complexity, making them an excellent tool for tasks that require large amounts of data to be processed in a more manageable form.

**K-Means:**

K-Means is a popular unsupervised machine learning algorithm used for clustering, where the goal is to partition a dataset into groups (clusters) based on similarity. The algorithm works by initializing K centroids, which represent the center of each cluster. It then assigns each data point to the nearest centroid, forming K clusters. Afterward, the centroids are recalculated as the mean of the data points assigned to each cluster, and the assignment and centroid update steps are repeated iteratively until convergence. The main objective is to minimize the within-cluster variance, or the sum of squared distances between data points and their corresponding

centroids. K-Means is widely used for tasks like customer segmentation, image compression, and anomaly detection. The process is illustrated in fig **3.8**.
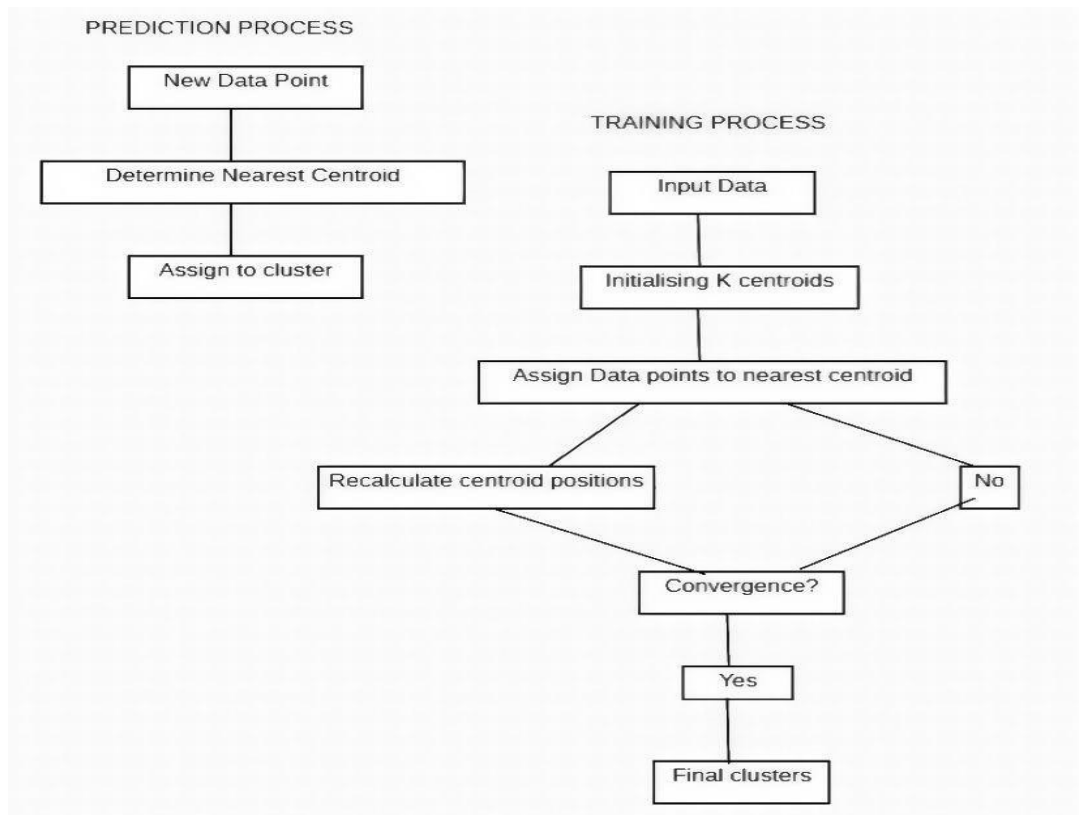


**Fig. 3.8: Kmeans**

One of the key strengths of K-Means is its simplicity and efficiency, especially when dealing with large datasets. However, it does have some limitations, such as requiring the user to specify the number of clusters (K) in advance, which can be difficult when the optimal number of clusters is unknown. Additionally, K-Means assumes that clusters are spherical and evenly sized, making it less effective for datasets with clusters of varying shapes or densities. Despite these limitations, K-Means remains a robust and scalable algorithm for clustering tasks, often serving as a baseline for more advanced clustering methods.

**DBSCAN:**

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a powerful clustering algorithm that can discover clusters of arbitrary shape and is particularly adept at handling noise in the data. Unlike K-Means, which requires the user to specify the number of clusters, DBSCAN works by identifying regions of

19

high data point density and expanding clusters from these dense regions. It does this by defining two main parameters: epsilon ($\varepsilon$), which determines the maximum distance between two points for them to be considered neighbors, and minPts, which sets the minimum number of points required to form a dense region. Points that do not meet these criteria are considered outliers or noise. This density-based approach allows DBSCAN to naturally identify clusters of varying shapes, unlike K-Means, which assumes spherical clusters.

DBSCAN is particularly effective in situations where clusters have irregular shapes or where there are significant amounts of noise. It does not require the number of clusters to be specified beforehand, which is a significant advantage over algorithms like K-Means. However, DBSCAN's performance is sensitive to the choice of the $\varepsilon$ and minPts parameters, and choosing appropriate values for these parameters can be challenging. Additionally, DBSCAN may struggle with datasets that contain clusters of varying densities, as it may classify points from denser clusters as noise if they are not sufficiently close to a dense region.

## DEC (Deep Embedded Clustering):

Deep Embedded Clustering (DEC) is a model that combines deep learning and clustering to improve the quality of clustering by learning a low-dimensional, feature-rich representation of the data that is more suitable for clustering. DEC integrates an autoencoder-like network architecture with clustering objectives as shown in fig **3.9**. The deep neural network learns an embedded representation of the input data that captures its underlying structure, and simultaneously, the clustering is performed by assigning data points to clusters in this low-dimensional space. The model minimizes a combined loss function, which includes both the reconstruction loss (like in autoencoders) and a clustering loss that encourages the network to group similar data points together in the latent space.

The key advantage of DEC over traditional clustering methods, such as K-Means, is its ability to leverage deep learning for feature extraction, enabling it to handle high-dimensional and complex data. DEC can automatically learn meaningful features from the data, which may not be possible using traditional clustering techniques that rely solely on distance metrics in the raw feature space. This makes DEC particularly

useful for clustering tasks involving high-dimensional data, such as image clustering or text clustering. Additionally, DEC can produce more accurate clusters by using a more sophisticated representation of the data, rather than relying on predefined distance metrics. However, DEC can be computationally expensive due to the use of deep learning models and may require careful tuning of the network architecture and clustering parameters.
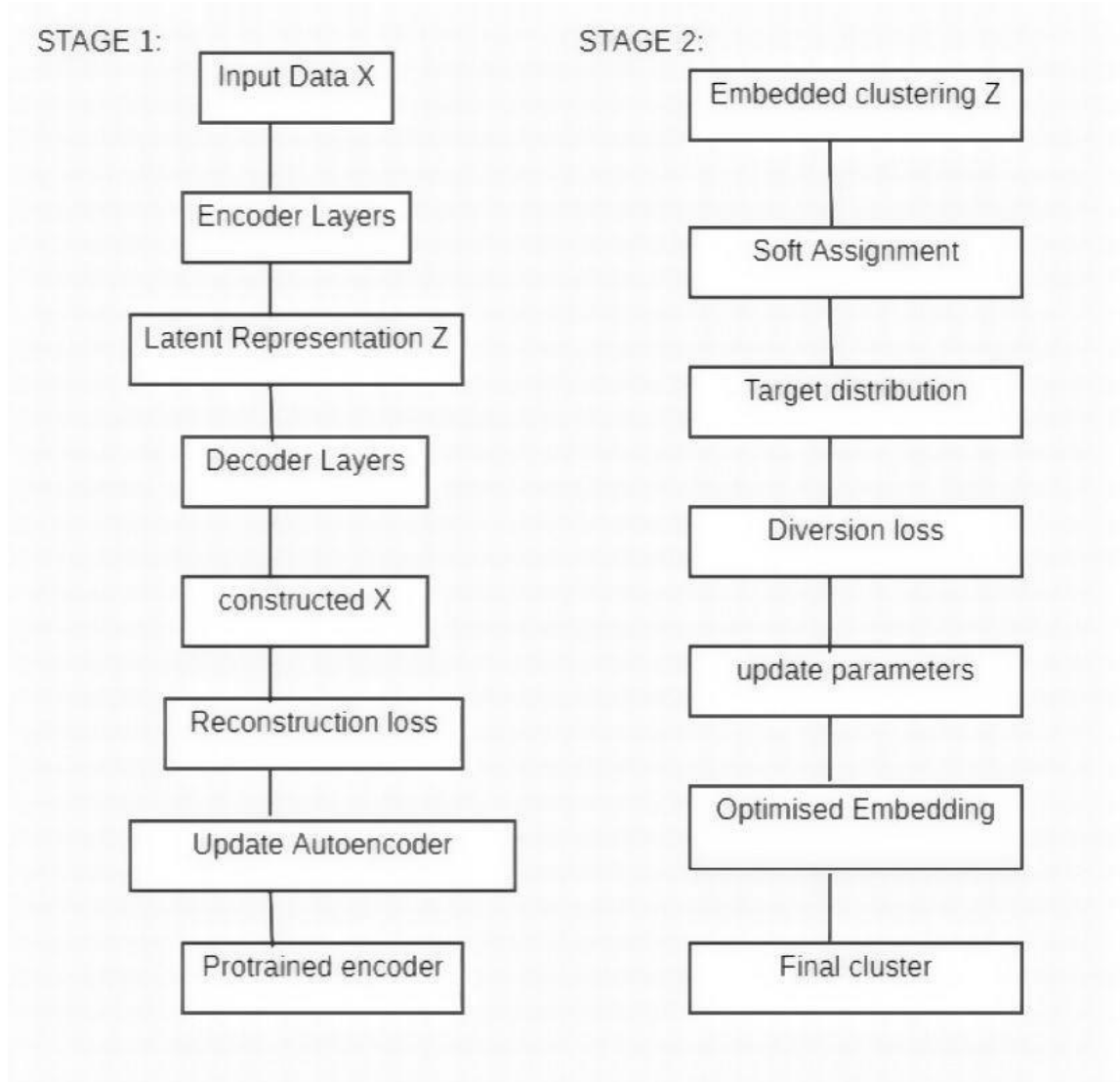
STAGE 1:

- Input Data X
- Encoder Layers
- Latent Representation Z
- Decoder Layers
- constructed X
- Reconstruction loss
- Update Autoencoder
- Protrained encoder

STAGE 2:

- Embedded clustering Z
- Soft Assignment
- Target distribution
- Diversion loss
- update parameters
- Optimised Embedding
- Final cluster

**Fig. 3.9: DEC**

**Benefits of DEC:**

- Dynamic Adaptation: DEC adapts its latent space representation dynamically as the clustering process evolves, ensuring that the learned features remain relevant and improve clustering accuracy.

- Handling Complex Data: DEC is highly effective for non-linear and high-

21

dimensional data, such as images, text, and audio, where traditional distance-based clustering methods often fail.

- Unsupervised Learning: As an unsupervised model, DEC does not require labeled data, making it suitable for exploratory data analysis in domains where labels are scarce or expensive to obtain.

**Applications:**

DEC is widely used in scenarios where traditional clustering techniques fall short due to the complexity of data, including:

- **Image Clustering:** Grouping images based on visual similarity or latent features.
- **Text Clustering:** Organizing text documents into topics or themes.
- **Genomic Data Analysis:** Identifying patterns and groups in high-dimensional genetic data.
- **Anomaly Detection:** Distinguishing abnormal data points in fraud detection or cybersecurity.

By combining the strengths of deep learning and clustering, DEC offers a powerful solution for unsupervised learning tasks, enabling data scientists to uncover patterns and insights in complex datasets effectively.

**Standard Scaler:**

StandardScaler is a preprocessing technique that standardizes the dataset by removing the mean and scaling it to unit variance. This scaling transforms the features to follow a standard normal distribution, with a mean of $0$ and a standard deviation of $1$. The transformation ensures that all features contribute equally to the model, especially when their magnitudes differ significantly.

However, in the presence of outliers, StandardScaler may not guarantee balanced feature scales because the outliers significantly influence the empirical mean and standard deviation. These extreme values can stretch or shrink the feature range, leading to distorted results.

The formula used by StandardScaler is:

$$Z = \frac{x - \mu}{\sigma}$$

where: x is the feature value,
$\quad\quad\mu$ is the mean of the feature,

σ is the standard deviation of the feature.

While StandardScaler works well for datasets without significant outliers, its performance can degrade if outliers are present. For such cases, combining StandardScaler with RobustScaler or removing outliers beforehand can produce better results.

**Label encoding:**

In machine learning projects, we usually deal with datasets having different categorical columns where some columns have their elements in the ordinal variable category for e.g a column income level having elements as low, medium, or high in this case we can replace these elements with 1,2,3. where 1 represents '*low'* 2 '*medium'* and 3' *high'.* Through this type of encoding, we try to preserve the meaning of the element where higher weights are assigned to the elements having higher priority.

Label Encoding is a technique that is used to convert categorical columns into numerical ones so that they can be fitted by machine learning models which only take numerical data. It is an important pre-processing step in a machine-learning project.

## 3.2.6 PERFORMANCE METRICS

**Clustering Accuracy**:

Measures how well the clustering algorithm groups similar data points together. It evaluates the degree of alignment between the predicted clusters and the ground truth labels, if available.

**Silhouette Score**:

A metric used to assess the quality of the clusters. It considers both cohesion (how close points in a cluster are) and separation (how distinct a cluster is from others). A higher silhouette score indicates better-defined clusters.

$$S(i) = \frac{(b - a)}{\max(a, b)} \qquad\qquad \text{(Eq- 3.3)}$$

where **a (Intra-cluster distance):** This measures how "close" a data point is to other

data points within its own cluster. A lower value indicates that the data point is more tightly clustered with its group.

**b (Nearest cluster distance):** This measures how "far" a data point is from other data points in the next closest cluster. A higher value indicates that the data point is well-separated from other clusters.

**Inertia (for K-Means)**:

Inertia measures the sum of squared distances between data points and their corresponding cluster centroids.

**Low Inertia**: Indicates that data points are closely packed around their respective centroids, reflecting well-separated clusters.

**High Inertia**: Suggests that data points are widely spread, possibly due to poor clustering or inappropriate choice of k.

**Adjusted Rand Index (ARI)**:

A metric that compares the similarity between two clusterings (the predicted clustering and the ground truth), adjusting for chance. ARI ranges from -1 to 1, with 1 indicating perfect clustering.

**DBSCAN Homogeneity**:

Homogeneity measures the degree to which clusters generated by DBSCAN (or any clustering algorithm) contain points belonging to the same class (ground truth label).

It evaluates whether each cluster includes only data points from a single class.

**Mean Squared Error (MSE)** (for Autoencoder):

Used to measure the difference between the input data and its reconstruction. Lower MSE indicates better reconstruction and feature extraction by the Autoencoder.

R-squared ($R^2$) gives an indication of how well the model fits the data by measuring the percentage of the variance in the dependent variable (Yield) that can be

predicted from the independent variables (features). The Mean Squared Error (MSE) between the expected and actualvalues is a measurement. Larger errors are given more weight when calculating the average of the squared discrepancies between the actual and anticipated values MSE is calculated using the formula:

$$MSE = (1/n) \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 \qquad\text{———————— (Eq-3.4)}$$

Where: n is the number of observations in the dataset.

   yi is the actual value of the observation.

   Y^i is the predicted value of the ith observation.

**Mean absolute error:**

Mean Absolute Error (MAE) is a measure of the average size of the mistakes in a collection of predictions, without taking their direction into account. It is measured as the average absolute difference between the predicted values and the actual values and is used to assess the effectiveness of a regression model.

The MAE loss function formula:

$$MAE = (1/n) \sum_{i=1}^{n} |y_i - \hat{y}_i| \qquad\text{———————— (Eq-3.5)}$$

where: yi :Actual value for the ith observation
   y^i : Calculated value for the ith observation
   n : Total number of observations

**Cluster Purity (for DEC):**

**Cluster Purity** measures the extent to which each cluster contains data points primarily from one ground truth class. It evaluates how well the clustering algorithm groups data points from the same class into the same cluster. Higher purity implies that the model has effectively grouped data points from the same class together.

**Outlier Detection Rate (for DBSCAN):**

Measures how effectively DBSCAN identifies outliers or noise points. A higher detection rate indicates better handling of noise and irrelevant data points in the clustering process.

# CHAPTER – 4

## IMPLEMENTATION

The implementation of this project involves a series of steps that focus on clustering techniques and deep learning for analyzing and classifying data. The first stage is data collection and preparation, where raw data is gathered and preprocessed to ensure consistency, missing value handling, and feature scaling. This is followed by implementing various clustering algorithms such as K-Means, DBSCAN, and DEC. K-Means is employed for partitioning data into predefined clusters based on distance from centroids, while DBSCAN is used to identify clusters of varying shapes and densities, also distinguishing outliers. DEC (Deep Embedded Clustering) is utilized to enhance clustering results by incorporating deep learning, where the model simultaneously learns a representation of the data while optimizing for cluster assignments. Autoencoders are also applied for feature extraction and dimensionality reduction, compressing input data to a lower-dimensional space, capturing essential features for clustering. The project also involves evaluating the performance of each algorithm using metrics like Silhouette Score, Inertia, and Adjusted Rand Index to compare and determine the best-performing clustering approach for the given dataset. These models and techniques work together to provide a robust and scalable solution for unsupervised data analysis and classification.

## 4.1 TOOLS USED

**Python**

Python is an interpreted, object-oriented, high-level programming language with dynamic semantics developed by Guido van Rossum. It was originally released in 1991. Designed to be easy as well as fun, the name "Python" is a nod to the British comedy group Monty Python. Python is an interpreted, object-oriented, high-level programming language with dynamic semantics developed by Guido van Rossum. It was originally released in 1991. Designed to be easy as well as fun, the name "Python" is a nod to the British comedy group Monty Python. Python has a reputation as a beginner-friendly language, replacing Java as the most widely used introductory

language because it handles much of the complexity for the user, allowing beginners to focus on fully grasping programming concepts rather than minute details.

Python is used for server-side web development, software development, mathematics, and system scripting, and is popular for Rapid Application Development and as a scripting or glue language to tie existing components because of its high-level, built-in data structures, dynamic typing, and dynamic binding. Program maintenance costs are reduced with Python due to the easily learned syntax and emphasis on readability. Additionally, Python's support of modules and packages facilitates modular programs and reuse of code. Python is an open source community language, so numerous independent programmers are continually building libraries and functionality for it. Professionally, Python is great for backend web development, data analysis, artificial intelligence, and scientific computing. Developers also use Python to build productivity tools, games, and desktop apps.

**Google Colaboratory:**

Google Colab, short for Colaboratory, is a free cloud-based platform provided by Google that allows users to write and execute Python code collaboratively in a Jupyter Notebook environment. Google Collaboratory notebook, is designed to facilitate machine learning (ML) and data science tasks by providing a virtual environment, Google colab python with access to free GPU resources.

**Benefits of Google Colab :**

Google Colab offers several benefits that make it a popular choice among data scientists, researchers, and machine learning practitioners.

**Key features of Google Collaboratory notebook include:**
- **Free Access to GPUs:** Colab offers free GPU access, which is particularly useful for training machine learning models that require significant computational power.
- **No Setup Required**: Colab runs in the cloud, eliminating the need for users to set up and configure their own development environment. This makes it convenient for quick coding and collaboration.

- **Collaborative Editing:** Multiple users can work on the same Colab notebook simultaneously, making it a useful tool for collaborative projects.
- **Integration with Google Drive**: Colab is integrated with Google Drive, allowing users to save their work directly to their Google Drive account. This enables easy sharing and access to notebooks from different devices.
- **Support for Popular Libraries**: Colab comes pre-installed with many popular Python libraries for machine learning, data analysis, and visualization, such as TensorFlow, PyTorch, Matplotlib, and more.
- **Easy Sharing**: Colab notebooks can be easily shared just like Google Docs or Sheets. Users can provide a link to the notebook, and others can view or edit the code in real- time.

**Clustering Libraries:**

For implementing clustering algorithms in this project, several libraries have been used, including Scikit-learn, TensorFlow, and Keras. Scikit-learn is one of the most popular machine learning libraries in Python, which provides a comprehensive suite of tools for clustering, classification, regression, and dimensionality reduction. It supports algorithms like K-Means, DBSCAN, and various evaluation metrics for model validation. TensorFlow, coupled with Keras, is leveraged for more complex models like Deep Embedded Clustering (DEC), which integrates deep learning into the clustering process. These libraries allow for the efficient implementation of unsupervised learning models, optimizing the process for handling large datasets and ensuring scalability across various clustering tasks. With their strong community support and integration with other tools, these libraries form the backbone of clustering workflows.

**Autoencoder (Keras):**

Autoencoders are implemented using Keras, a high-level neural networks API written in Python. Keras provides a simple and easy-to-use interface for building deep learning models, including autoencoders. Autoencoders are used in this project for dimensionality reduction, where the input data is compressed into a lower-dimensional representation and then reconstructed back to the original form. This helps in feature extraction, capturing the most important aspects of the data while

discarding irrelevant information. The Keras library simplifies building these models by offering various layers and activation functions, making it easy to experiment with different architectures and training strategies. With Autoencoders, the model learns to represent the data efficiently, improving the performance of subsequent clustering algorithms.

**Matplotlib and Seaborn (Visualization):**

Matplotlib and Seaborn are essential Python libraries for data visualization. Matplotlib is used for creating static, animated, and interactive plots in Python. It allows for extensive customization of plots, including adding titles, labels, and legends. It supports various output formats such as PNG, PDF, and SVG, making it versatile for different use cases. Seaborn, built on top of Matplotlib, enhances its functionality by providing high-level interfaces for drawing attractive statistical graphics. It is particularly useful for visualizing complex relationships in the data, such as heatmaps, correlation matrices, and distribution plots. Seaborn's default styling options and color palettes allow for quick creation of aesthetically pleasing and informative visualizations. These tools are crucial for exploring the data and evaluating the results of clustering models in this project.

**TensorFlow (Deep Learning):**

TensorFlow is an open-source machine learning library developed by Google, used for building and training deep learning models. TensorFlow is highly efficient and supports a wide range of machine learning tasks, from simple linear regressions to complex deep neural networks. In this project, TensorFlow is used for the implementation of DEC (Deep Embedded Clustering), which involves deep learning for clustering tasks. TensorFlow's ability to handle large datasets and perform parallel computations on GPUs makes it suitable for training models on high-dimensional data. With Keras integrated into TensorFlow, building neural network architectures becomes significantly easier, allowing for seamless experimentation with deep learning models that enhance clustering outcomes by learning data representations.

**Scikit-learn (Clustering Algorithms):**

Scikit-learn is a robust and widely-used Python library for machine learning,

offering a simple interface to implement clustering algorithms like K-Means, DBSCAN, and others. It is a versatile tool for both supervised and unsupervised learning, providing easy-to-use APIs for model fitting, predictions, and performance evaluation. Scikit- learn's implementation of K-Means and DBSCAN allows for rapid deployment and evaluation of clustering tasks on large datasets, making it a key component of this project. Scikit-learn also includes various utilities for data preprocessing, cross- validation, and metric evaluation, making it an all-encompassing library for machine learning tasks.

**Pandas (Data Manipulation):**

Pandas is a Python library used for data manipulation and analysis, providing efficient data structures like DataFrames that allow for easy handling of structured data. In this project, Pandas is essential for reading, cleaning, and preprocessing the data. It supports various file formats like CSV, Excel, and SQL databases, making it easy to load and explore large datasets. Pandas allows for fast operations like filtering, merging, and reshaping datasets, which are necessary steps in preparing the data for clustering. The library's integration with other libraries like Matplotlib and Seaborn further simplifies the process of visualizing the data and the results of clustering models.

## 4.2 CODE

*!pip install xlrd==1.2.0*
*!pip install squarify*
*!pip install -q kaggle*
*#!/bin/bash*
**!kaggle datasets download olistbr/brazilian-ecommerce**
*!unzip brazilian-ecommerce.zip -d /content/data/*

*import numpy as np*
*import pandas as pd*
*import datetime as dt*
*import matplotlib.pyplot as plt*
*import seaborn as sns import squarify*

**# Reading data from all tables in the dataset**

```
customers_df = pd.read_csv('./data/olist_customers_dataset.csv')

items_df = pd.read_csv('./data/olist_order_items_dataset.csv')

products_df= pd.read_csv('./data/olist_products_dataset.csv')

category_translation_df=pd.read_csv('./data/product_category_name_translation.csv')

orders_df = pd.read_csv('./data/olist_orders_dataset.csv')

payments_df = pd.read_csv('./data/olist_order_payments_dataset.csv')


# Merging all dataframes

merged_df = orders_df.merge(items_df, on='order_id')

merged_df = merged_df.merge(payments_df, on='order_id')

merged_df = merged_df.merge(products_df, on='product_id')

merged_df = merged_df.merge(customers_df, on='customer_id')

merged_df=merged_df.merge(category_translation_df, on='product_category_name')

merged_df.head()

 merged_df.info()
#convert column that relate to date to 1 type - datetime dtype

date_feature=['order_purchase_timestamp','order_approved_at','order_delivered_carrier_date','order_delivered_customer_date','order_estimated_delivery_date','shipping_limit_date' ]
for i in date_feature:
    merged_df[i]=pd.to_datetime(merged_df[i],errors ='raise', utc = False)


rfm_d=merged_df[['customer_unique_id','order_id','product_category_name_english','price','freight_value','payment_type','order_status','order_purchase_timestamp','customer_state','payment_value']]
rfm_df.info()
 rfm_df=rfm_df[rfm_df['order_status'] =='delivered']


rfm_df.info()
#EDA

 import matplotlib.pyplot as plt
```

31

```python
import numpy as np
import pandas as pd
import pydotplus as pdp
import scipy.stats as st
import seaborn as sbn
import plotly.express as px
from sklearn.decomposition import PCA
import warnings
warnings.filterwarnings('ignore')
from IPython.display import Image
from sklearn import linear_model, svm, tree
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, mean_squared_error, precision_score, f1_score, recall_score#, plot_confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import BernoulliNB, GaussianNB from sklearn.neighbors import KNeighborsClassifier
from sklearn.tree import DecisionTreeClassifier
from scipy.cluster import hierarchy from scipy.spatial.distance import cdist
from sklearn.cluster import KMeans, AgglomerativeClustering
from sklearn.preprocessing import MinMaxScaler, StandardScaler
import statsmodels.api as sm
from statsmodels.formula.api import ols
from statsmodels.stats.multicomp import pairwise_tukeyhsd
import scipy.stats as stats


prod_gb1=rfm_df.groupby(by=['order_id'])['price'].count().reset_index()
prod_gb1.columns=['order_id','ProdQuantity']
display(prod_gb1.sort_values(by = 'Prod Quantity', ascending = False))


prod_gb2=pd.DataFrame(rfm_df.groupby(["order_id"],as_index=False).agg({"price":"nunique"}))
prod_gb2.columns=['order_id','UniqueProd']
display(prod_gb2.sort_values(by='UniqueProd',ascending=False))
```

```python
result1=pd.merge(prod_gb1, prod_gb2, on = 'order_id') display(result1)
display(result1.sort_values(by = 'Prod Quantity', ascending = False).head(10))


prod_tb1=pd.DataFrame(index=sorted(result1['ProdQuantity'].unique().tolist()))
prod_tb1['NumofProd'] = prod_tb1.index
prod_gb11 =result1.groupby('ProdQuantity')['ProdQuantity'].count()
prod_tb1['Amount'] = prod_gb11
display(prod_tb1)


prod_tb2 = pd.DataFrame(index=sorted(result1['Unique Prod'].unique().tolist()))
prod_tb2['NumofProd'] = prod_tb2.index
prod_gb11 =result1.groupby('Unique Prod')['Unique Prod'].count()
prod_tb2['Amount'] = prod_gb11
display(prod_tb2)


ord_gb2 = pd.DataFrame(rfm_df.groupby(["customer_unique_id"], as_index =
False).agg({"order_id":"nunique"}))
ord_tb1 = pd.DataFrame(index=sorted(ord_gb2.order_id.unique().tolist()))
ord_gb2= ord_gb2.groupby('order_id')['order_id'].count()
ord_tb1['Amount'] = ord_gb2 display(ord_tb1.sort_values(by = 'Amount',
ascending = False))
ord_unique = rfm_df.drop_duplicates(subset=['order_id'], keep='first')
display(ord_unique)
time_tb=
ord_unique.groupby(ord_unique['order_purchase_timestamp'].dt.strftime('%y-
%m')).agg (total_orders = ('order_id' , 'count'))
time_tb.total_orders.sum()display(time_tb)


import plotly.graph_objects as go
fig=go.Figure(data=go.Scatter(x=time_tb.index,y=time_tb["total_orders"],mode
='lines+markers'))
fig.update_xaxes(title_text="timestamp(yy-mm)")
fig.update_yaxes(title_text="Number of Orders")
```

33

```
fig.update_layout(title_text='Time stamp vs Number of Orders', title_x=0.5)
fig.update_traces(textposition="bottom right") fig.show()

cate_gb=rfm_df.groupby('product_category_name_english').aggregate({'payment_value':
['count', 'mean']})

cate_gb.info()
#cate_gb['category'] = cate_gb.index
display(cate_gb)
cate_gb.columns = ['_'.join(x) if x[1] else x[0] for x in cate_gb.columns]
cate_gb = cate_gb.reset_index()
cate_gb = cate_gb.round({'payment_value_mean':2})

fig=  px.scatter(cate_gb,  x="payment_value_mean",  y="payment_value_count",
size="payment_value_count",range_color="Viridis",hover_name="product_categor
y_name_english",size_max=60)


fig.update_xaxes(title_text="Means of Orders  Value")
fig.update_yaxes(title_text="Quantity of Sold Products")
fig.update_layout(title_text='Means of order value vs Quantity of sold products',
title_x=0.5)
fig.show()

pay = rfm_df[['payment_type', 'payment_value','order_id']].drop_duplicates()

display(pay)

payme_gb = pay.groupby(by = ['payment_type',
'payment_value'])['order_id'].count().reset_index()

display(payme_gb)

fig=px.scatter(payme_gb,x="payment_value",y="order_id",color="payment_type",
#marginal_x="box", marginal_y="violin", size = 'order_id', size_max = 60)
fig.update_xaxes(title_text="Value of Orders")
fig.update_yaxes(title_text="Number of Orders")
fig.update_layout(title_text='Value of Orders vs Number of Orders', title_x=0.5)
fig.show()
```

# RFM analysis

```
present_day=rfm_df['order_purchase_timestamp'].max()+dt.timedelta(days=2)
print("Latest date in dataset: ", rfm_df['order_purchase_timestamp'].max()) rfm_df
rfm_df=rfm_df[['customer_unique_id','order_id','payment_value','order_purchase_ti
mestamp']]
rfm_df.head()
rfm_df=rfm_df.groupby('customer_unique_id').agg({'order_id':'nunique','order_purc
hase_ti mestamp':['min','max'],'payment_ value':'sum'}).reset_index()
rfm_df.columns=["customer_unique_id","frequency","first_order_date","last_orde
r_date","monetary "]
```

#Obtain recency value for each customer adding new column

```
rfm_df['recency']= rfm_df['last_order_date'].apply(lambda x: (present_day - x).days)
```

#Deleting date column because it's not necessary

```
rfm_df=rfm_df[['customer_unique_id','recency','frequency','monetary']]
 #Show results
rfm_df.head()
rfm_df = rfm_df.sample(n=10000, random_state=42) rfm_df.info()
rfm_df.describe() rfm_df_outlier= rfm_df.copy()
rfm_df.to_csv('rfm_dataframe.csv')
```

# Identify Outliers

```
rfm_iqr = rfm_df_outlier.copy()
rfm_iqr
 attributes = ['monetary','frequency','recency'] plt.rcParams['figure.figsize'] =
[10,8]
sns.boxplot(data=rfm_iqr[attributes],orient="v",palette="Set2",whis=1.5,saturation
=1, width=0.7)
plt.title("Outliers Variable Distribution", fontsize=14,fontweight='bold')
plt.ylabel("Range",fontweight='bold')
plt.xlabel("Attributes", fontweight = 'bold')
```

**IQR**

# *Note for normal distribution*

*from pandas.core.reshape.encoding import DataFrame*

# *Removing (statistical) outliers for monetary*
*rfm_features=['monetary']*
 *def iqr_outlier (df:DataFrame,ft:list):*

    *for i in ft:*

       *Q1 = df[i].quantile(.25)*

       *Q3 = df[i].quantile(.75)*

       *IQR = Q3 - Q1*

    *df = df[(df[i] < Q1 - 1.5\*IQR) | (df[i] > Q3 + 1.5\*IQR)]*

 *return df rfm_i = pd.DataFrame()*

*rfm_i['monetary']=iqr_outlier(rfm_df_outlier,rfm_features)['monetary']*
*display(rfm_i)*

# ***Z-score***

*rfm_zs= pd.DataFrame#rfm_df_outlier.copy()*

*rfm_features=['monetary'] def zscore_outliers(df:DataFrame,ft:list,*

*threshold):*

  *for col in ft:*

    *z_scores = (df[col] - df[col].mean()) / df[col].std()*

  *# Identify outliers based on threshold*

  *outliers = pd.DataFrame(df[col][abs(z_scores) >*

 *threshold]) return outliers*

*rfm_zs=zscore_outliers(rfm_df_outlier,rfm_features,3)  display(rfm_zs)*

# ***Filter outliers***

*iqr = rfm_i.index.tolist()*

*zs = rfm_zs.index.tolist()*

*outliers = list(set(iqr) & set(zs)) display(len(outliers))*

## ## Rescaling the Attributes

*It is extremely important to rescale the variables so that they have a comparable scale.| There are two common ways of rescaling:*
*1. Min-Max scaling*
*2. Standardisation (mean-0, sigma-1)*
*Here, we will use Standardisation Scaling.*

*from sklearn.preprocessing import StandardScaler*
## # Rescaling the attributes
*rfm = rfm_df[['monetary','frequency', 'recency']]*

## # Instantiate

*scaler = StandardScaler()*

## # fit_transform

*rfm_scaled = scaler.fit_transform(rfm) rfm_scaled.shape*
*rfm_df_scaled = pd.DataFrame(rfm_scaled)*
*rfm_df_scaled.columns = ['monetary', 'frequency', 'recency']*
*rfm_df_scaled.head()*
*from sklearn.metrics import silhouette_samples import matplotlib.cm as cm*
*from yellowbrick.cluster import KElbowVisualizer*
*import sklearn*
*from sklearn.preprocessing import StandardScaler*
*from sklearn.cluster import KMeans*
*from sklearn.metrics import silhouette_score*
*from scipy.cluster.hierarchy import linkage*
from scipy.cluster.hierarchy import dendrogram
*from scipy.cluster.hierarchy import cut_tree*

## Finding the Optimal Number of Clusters

#### Elbow Curve to get the right number of Clusters

*A fundamental step for any unsupervised algorithm is to determine the optimal number of clusters into which the data may be clustered. The Elbow Method is one of the most popular methods to determine this optimal value of k.*

# Elbow-curve/SSD

*ssd = [] range_n_clusters = [2, 3, 4, 5, 6, 7, 8]*
*for num_clusters in range_n_clusters:*
*kmeans=KMeans(n_clusters=num_clusters,max_iter=50)*
*kmeans.fit(rfm_df_scaled)*
*ssd.append(kmeans.inertia_)*

*# plot the SSDs for each n_cluster*
*plt.plot(range_n_clusters, ssd)*
*plt.plot Distorition score elbow để xác định số lượng cụm tối ưu*
*rfmk_means = KMeans()*

*elbow = KElbowVisualizer(rfmk_means, k=(2, 8))*
*elbow.fit(rfm_df_scaled) elbow.show()*

### Silhouette Analysis

*{silhouette score}={p-q}/{max(p,q)}*

*p is the mean distance to the points in the nearest cluster that the data point is not a part of.*

*q is the mean intra-cluster distance to all the points in its own cluster.*

- *The value of the silhouette score range lies between -1 to 1.*

- *A score closer to 1 indicates that the data point is very similar to other data points in the cluster.*

- *A score closer to -1 indicates that the data point is not similar to the data points in its cluster.*

# Silhouette analysis

```
range_n_clusters = [2, 3, 4, 5, 6, 7, 8]
for num_clusters in range_n_clusters:
    # intialise kmeans
    kmeans=KMeans(n_clusters=num_clusters,max_iter=50)
kmeans.fit(rfm_df_scaled)
cluster_labels = kmeans.labels_

# silhouette score
silhouette_avg = silhouette_score(rfm_df_scaled, cluster_labels)
print("For n_clusters={0}, the silhouette score is {1}".format(num_clusters,
silhouette_avg))
```

## KMeans

```
# Final model with k=4
kmeans = KMeans(n_clusters=4, max_iter=50)
kmeans.fit(rfm_df_scaled)
kmeans.labels_
# assign the label
rfm_df['clus_kmeans'] = kmeans.labels_ rfm_df.head()

rfm_df['clus_kmeans'].value_counts().plot(kind='bar')

plt.xlabel('Cluster')

 plt.ylabel('Amount of person')

plt.title('Amount of person per cluster')

plt.show()

# Box plot to visualize Cluster Id vs monetary

sns.boxplot(x='clus_kmeans', y='monetary', data=rfm_df)
```

**# Box plot to visualize Cluster Id vs Frequency**

```
sns.boxplot(x='clus_kmeans', y='frequency', data=rfm_df)

#sns.jointplot(x='clus_kmeans',y='frequency',data=rfm_df,kind  ='hist',color

='red') # Box plot to visualize Cluster Id vs Recency
```

*sns.boxplot(x='clus_kmeans', y='recency', data=rfm_df)*

## Hierarchical Clustering

*Hierarchical clustering involves creating clusters that have a predetermined ordering from top to bottom. For example, all files and folders on the hard disk are organized in a hierarchy.*

*There are two types of hierarchical clustering,*

*- Divisive*

*- Agglomerative.*

*import sklearn*

*from sklearn.preprocessing import StandardScaler*

*from sklearn.cluster import KMeans*

*from sklearn.metrics import silhouette_score*

*from scipy.cluster.hierarchy import linkage*

*from scipy.cluster.hierarchy import dendrogram*

*from scipy.cluster.hierarchy import cut_tree*

*# Average linkage*

*mergings = linkage(rfm_df_scaled, method="average", metric='euclidean')*

*dendrogram(mergings)*

*plt.show()*

*#### Cutting the Dendrogram based on K*
*# 4 clusters*
*cluster_labels = cut_tree(mergings, n_clusters=4).reshape(-1, )*
*cluster_labels*

*# Assign cluster labels*

*rfm_df['clus_hac'] = cluster_labels rfm_df.head()*

*# Plot Cluster Id vs Amount*

*sns.boxplot(x='clus_hac', y='monetary', ata=rfm_df)*
*# Plot Cluster Id vs Frequency*

*sns.boxplot(x='clus_hac', y='frequency',*

*data=rfm_df)*

*# Plot Cluster Id vs Recency*

40

```
sns.boxplot(x='clus_hac', y='recency', data=rfm_df)
rfm_df.to_csv('rfmclus.csv', index=False)

import pandas as pd import numpy as np
from sklearn.preprocessing import StandardScaler
from tensorflow.keras.models import Sequential, Model
from tensorflow.keras.layers import Dense, Input
from sklearn.manifold import TSNE
import matplotlib.pyplot as plt

# Define Autoencoder model input_dim = rfm_df_scaled.shape[1]
encoding_dim = 2
input_layer = Input(shape=(input_dim,))

encoder = Dense(encoding_dim, activation='relu')(input_layer)
decoder =  Dense(input_dim, activation='sigmoid')(encoder)

autoencoder = Model(inputs=input_layer, outputs=decoder)
autoencoder.compile(optimizer='adam', loss='mse')
# Train Autoencoder model

autoencoder.fit(rfm_df_scaled, rfm_df_scaled, epochs=50, batch_size=32)
# Use Autoencoder to encode data
encoder = Model(inputs=input_layer, outputs=encoder)
encoded_data  = encoder.predict(rfm_df_scaled)

# Perform KMeans clustering
kmeans = KMeans(n_clusters=4, random_state=0).fit(encoded_data)

# Assign cluster labels to original data rfm_df['clus_autoencoder'] = kmeans.labels_
print(rfm_df)

# Plot Cluster Id vs Amount
sns.boxplot(x='clus_autoencoder', y='monetary',
data=rfm_df)

# Plot Cluster Id vs Frequency
sns.boxplot(x='clus_autoencoder', y='frequency',
```

*data=rfm_df)*

*# Plot Cluster Id vs Recency*

*sns.boxplot(x='clus_autoencoder', y='recency', data=rfm_df)*

*# Visualize clustering using t-SNE*

*tsne = TSNE(n_components=2, perplexity=30, random_state=0)*

*tsne_data = tsne.fit_transform(encoded_data)*

*plt.scatter(tsne_data[:, 0], tsne_data[:, 1])*

*plt.show()*

## Model Validation

*from sklearn.metrics import silhouette_score, calinski_harabasz_score*

*# Reconstruction Error*

*reconstruction_error = autoencoder.evaluate(rfm_df_scaled, rfm_df_scaled)*

# Silhouette Score

*silhouette = silhouette_score(encoded_data, kmeans.labels_) print(silhouette)*

# Calinski-Harabasz Index

*calinski_harabasz = calinski_harabasz_score(encoded_data, kmeans.labels_)*
*print(calinski_harabasz)*

##Omit outliers

*outlier_rm = rfm_df_outlier.drop(outliers)*

*outlier_rm*

## Rescaling the Attributes

*from sklearn.preprocessing import StandardScaler*

*# Rescaling the attributes*

rfm_df_outlier = outlier_rm[['monetary', 'frequency', 'recency']]

*# Instantiate*

*scaler = StandardScaler()*

*# fit_transform*

*rfm_scaled_outlier = scaler.fit_transform(rfm_df_outlier) rfm_scaled_outlier.shape*

*rfm_df_outlier_scaled = pd.DataFrame(rfm_scaled_outlier)*

```
rfm_df_outlier_scaled.columns = ['monetary', 'frequency', 'recency']
rfm_df_outlier_scaled.head()
```

## Finding the Optimal Number of Clusters

```
import sklearn
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
from scipy.cluster.hierarchy import linkage
from scipy.cluster.hierarchy import
dendrogram from scipy.cluster.hierarchy
import cut_tree

 # Elbow-curve/SSD
ssd = [] range_n_clusters = [2, 3, 4, 5, 6, 7, 8]
for num_clusters in range_n_clusters:
    kmeans = KMeans(n_clusters=num_clusters, max_iter=50)
kmeans.fit(rfm_df_outlier_scaled)
ssd.append(kmeans.inertia_)

# plot the SSDs for each n_clusters plt.plot(range_n_clusters, ssd)
plt.plot
from yellowbrick.cluster import KElbowVisualizer rfmk_means = KMeans()
elbow = KElbowVisualizer(rfmk_means, k=(2, 8))
elbow.fit(rfm_df_scaled)
elbow.show()

# Silhouette analysis
range_n_clusters = [2, 3, 4, 5, 6, 7,
8] for num_clusters in
range_n_clusters:
   # intialise kmeans
   Kmeans = KMeans(n_clusters=num_clusters, max_iter=50)
kmeans.fit(rfm_df_outlier_scaled)
cluster_labels = kmeans.labels_
```

43

```
# silhouette score
silhouette_avg = silhouette_score(rfm_df_outlier_scaled, cluster_labels)
print("For n_clusters={0}, the silhouette score is {1}".format(num_clusters,
silhouette_avg))

##KMeans
# Final model with k=4
kmeans = KMeans(n_clusters=4, max_iter=50) kmeans.fit(rfm_df_outlier_scaled)
kmeans.labels_

# assign the label
outlier_rm['clus_kmeans'] = kmeans.labels_ outlier_rm
outlier_rm['clus_kmeans'].value_counts().plot(kind='bar')
plt.xlabel('Cluster')
plt.ylabel('Amount of person')
plt.title('Amount of person per cluster')
plt.show()

## Hierarchical Clustering
# Average linkage
mergings = linkage(rfm_df_outlier_scaled, method="average", metric='euclidean')
dendrogram(mergings)

plt.show()
#4 clusters
cluster_labels = cut_tree(mergings, n_clusters=4).reshape(-1, ) cluster_labels

# Assign cluster labels

outlier_rm ['clus_hac'] = cluster_labels outlier_rm .head()

## Autoencoder

# Define Autoencoder model

input_dim_outlier = rfm_df_outlier_scaled.shape[1]
encoding_dim_outlier = 2
input_layer_outlier = Input(shape=(input_dim_outlier,))
```

```
encoder_outlier = Dense(encoding_dim_outlier,
ctivation='relu')(input_layer_outlier) decoder_outlier= Dense(input_dim_outlier,
activation='sigmoid')(encoder_outlier)
autoencoder_outlier = Model(inputs=input_layer_outlier, outputs=decoder_outlier)
autoencoder_outlier.compile(optimizer='adam', loss='mse')

# Train Autoencoder model
autoencoder_outlier.fit(rfm_df_outlier_scaled, rfm_df_outlier_scaled, epochs=50,
batch_size=32)

# Use Autoencoder to encode data
encoder_outlier    =    Model(inputs=input_layer_outlier, outputs=encoder_outlier)
encoded_data_outlier = encoder.predict(rfm_df_outlier_scaled)
outlier_rm

# Perform KMeans clustering
kmeans = KMeans(n_clusters=4, random_state=0).fit(encoded_data_outlier)

# Assign   cluster   labels   to   original   data
outlier_rm['clus_autoencoder'] = kmeans.labels_
print(outlier_rm)

##Model Validation

# Calinski-Harabasz Index
calinski_harabasz_outlier = calinski_harabasz_score(encoded_data_outlier,
kmeans.labels_)
 print(calinski_harabasz_outlier)

# Reconstruction Error
reconstruction_error_outlier = autoencoder.evaluate(rfm_df_outlier_scaled,
rfm_df_outlier_scaled)
reconstruction_error_outlier

# Plot Cluster Id vs Amount
sns.boxplot(x='clus_autoencoder', y='monetary',
data=rfm_df)
# Plot Cluster Id vs Frequency
```

```
sns.boxplot(x='clus_autoencoder', y='frequency',
data=rfm_df)

# Plot Cluster Id vs Recency
sns.boxplot(x='clus_autoencoder', y='recency', data=rfm_df)

# Silhouette Score
silhouette_outlier = silhouette_score(encoded_data_outlier, kmeans.labels_)
print(silhouette_outlier)

#Evaluate 2 circumstances
arr=np.array([[reconstruction_error,silhouette,calinski_harabasz],[reconstruction_er
ror_outlier,silhouette_outlier,cali nski_harabasz_outlier]])
evaluate=pd.DataFrame(arr,columns=['reconstruction_error','silhoutte','calinski_har
abasz'],index=['outlier','no outlier'])

print(evaluate)

import numpy as np import pandas as pd import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
evaluater=pd.read_csv('rfm_dataframe.csv')
from sklearn.metrics import silhouette_score, calinski_harabasz_score,
davies_bouldin_score, adjusted_rand_score, fowlkes_mallows_score

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from  sklearn.preprocessing  import  StandardScaler
from  sklearn.cluster  import DBSCAN
 from sklearn.metrics import silhouette_score

# DBSCAN parameters eps = 0.5 min_samples = 15
# DBSCAN Clustering
Dbscan = DBSCAN(eps=eps, min_samples=min_samples)
rfm_df_scaled['clus_dbscan'] = dbscan.fit_predict(rfm_df_scaled)
```

*# Show the clustering results*

*print(rfm_df_scaled['clus_dbscan'].value_counts())*

*# Shows the number of points in each cluster*

*plt.figure(figsize=(8, 6))*

*# Plot the clusters*

*plt.scatter(rfm_df_scaled.iloc[:,0],rfm_df_scaled.iloc[:,1],*

*c=rfm_df_scaled['clus_dbscan'],cmap='viridis', s=50)*

*plt.title("DBSCAN Clustering")*

*plt.xlabel('Feature 1')*

*plt.ylabel('Feature 2')*

*plt.colorbar(label='Cluster')*

*plt.show()*

*# Evaluate with Silhouette Score (not applicable to noise points, so ignore them)*

*valid_clusters = rfm_df_scaled['clus_dbscan'] != -1*

*silhouette_avg = silhouette_score(rfm_df_scaled[valid_clusters],*

*rfm_df_scaled.loc[valid_clusters, 'clus_dbscan'])*

*print(f"Silhouette Score:{silhouette_avg}")*

***DEC***

*!pip install tensorflow matplotlib*

*from tensorflow.keras.models import Model*

*from tensorflow.keras.layers import Input, Dense, BatchNormalization, Dropout*

*from tensorflow.keras.optimizers import Adam from sklearn.cluster import KMeans*

*from sklearn.metrics import silhouette_score*

*import numpy as np*

*import tensorflow as tf*

*# Fix seeds for reproducibility*

*np.random.seed(42)*

*tf.random.set_seed(42)*

*# Define model*

*input_dim = rfm_df_scaled.shape[1]*

*encoding_dim = 10*

47

```python
# Higher latent space dimensionality
input_layer=Input(shape=(input_dim,))
encoded = Dense(64, activation='relu')(input_layer)
encoded = BatchNormalization()(encoded)
encoded = Dropout(0.2)(encoded)
encoded = Dense(32, activation='relu')(encoded)
encoded = Dense(encoding_dim, activation='relu')(encoded)
decoded = Dense(32, activation='relu')(encoded)
decoded = Dense(64, activation='relu')(decoded)
decoded = Dense(input_dim, activation='sigmoid')(decoded)

# Compile and train autoencoder
autoencoder = Model(inputs=input_layer, outputs=decoded)
autoencoder.compile(optimizer=Adam(learning_rate=0.001),loss='binary_crossentropy')
autoencoder.fit(rfm_df_scaled,    rfm_df_scaled, epochs=200, batch_size=32,
verbose=1)

# Extract latent space
encoder_model = Model(inputs=input_layer, outputs=encoded)
encoded_data = encoder_model.predict(rfm_df_scaled)

# Perform clustering
Kmeans = KMeans(n_clusters=4, n_init=50, random_state=42)
kmeans.fit(encoded_data)
rfm_df['clus_dec'] = kmeans.labels_

# Evaluate clustering
silhouette_avg = silhouette_score(encoded_data, kmeans.labels_)
print(f"Silhouette Score for DEC: {silhouette_avg}")
```

## 4.3 RESULTS AND SCREENSHOTS

**Table 4.1:** Comparison Of Proposed vs Traditional Methods

| Method | Accuracy |
|---|---|
| K-Means | 58.2 |
| DBSCAN | 63.9 |
| Proposed Model | 66.4 |

The Table **4.1** illustrate the advantages of the proposed model over traditional clustering methods. While K-Means and DBSCAN rely solely on raw feature spaces and predefined distance metrics, the proposed model excels by leveraging deep learning to extract feature-rich representations, enabling more accurate and robust clustering, particularly in complex datasets.



**Fig 4.1**: Cluster distribution of individuals using K-Means based on RFM analysis.

The figure **4.1** illustrates the distribution of individuals across clusters derived using the K-Means algorithm, where the clustering is performed based on RFM

(Recency, Frequency, Monetary) analysis. The RFM framework is a well-known approach in customer segmentation that evaluates customers based on three criteria:

- **Recency (R):** How recently a customer made a purchase.
- **Frequency (F):** How often a customer makes purchases.
- **Monetary Value (M):** How much money a customer spends.



**Fig 4.2:** Dendrogram illustrating hierarchical clustering for RFM-based customer segmentation

The figure **4.2** represents a dendrogram, which is a tree-like diagram used to illustrate the arrangement of the clusters produced through hierarchical clustering based on RFM values.

The figure **4.3** illustrates the distribution of monetary values across the clusters identified using Deep Embedded Clustering (DEC). A box plot is used to visualize the spread, central tendency, and outliers for monetary values within each cluster.

**Fig 4.3:** Box plot showing the monetary value distribution for clusters obtained using Deep Embedded Clustering



**Fig 4.4:** Box plot showing the frequency distribution for clusters using Deep Embedded Clustering

51

The figure **4.4** illustrates the recency distribution across the clusters identified through Deep Embedded Clustering (DEC). Recency refers to how recently a customer made a purchase or interacted with the business, and this metric is used in customer segmentation to evaluate engagement patterns.
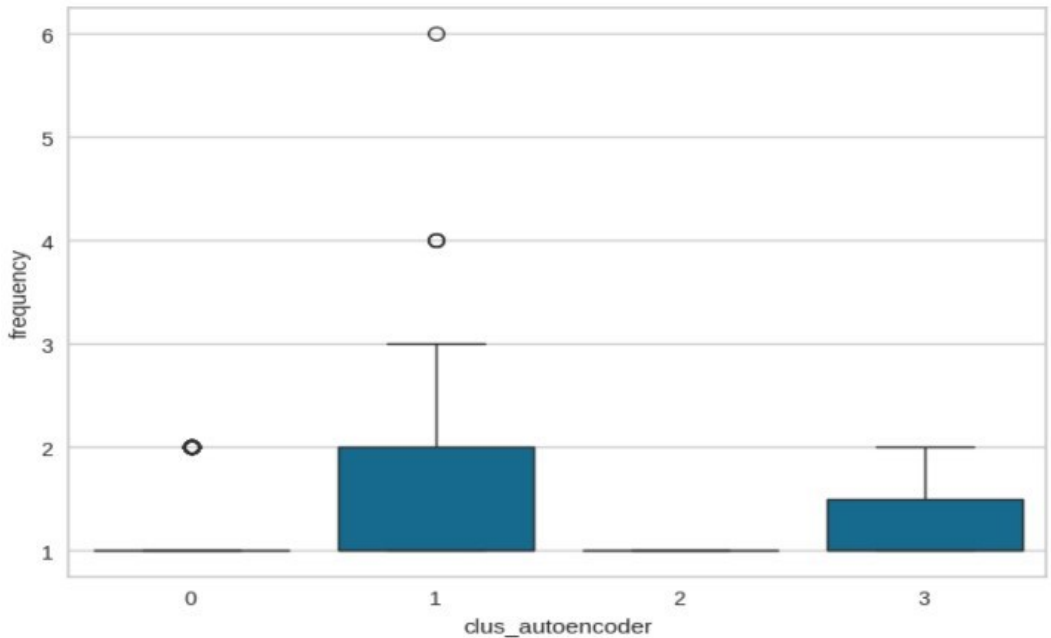


**Fig 4.5:** Box plot showing the recency distribution for clusters obtained using Deep Embedded Clustering

The box plot illustrates the frequency distribution for different customer clusters identified by DEC. In RFM analysis, frequency measures how often a customer makes purchases or interacts with the business during a given period.

The steady decline in loss demonstrates effective training progress in fig **4.6**. This model's outputs will likely serve as input features for clustering, enhancing the clustering algorithm's performance by providing meaningful representations of the original data.

```
Epoch 1/60
313/313 ─────────────────────── 3s 6ms/step - loss: 0.9253
Epoch 2/60
313/313 ─────────────────────── 2s 4ms/step - loss: 0.7062
Epoch 3/60
313/313 ─────────────────────── 2s 5ms/step - loss: 0.6673
Epoch 4/60
313/313 ─────────────────────── 1s 2ms/step - loss: 0.7276
Epoch 5/60
313/313 ─────────────────────── 2s 3ms/step - loss: 0.5488
Epoch 6/60
313/313 ─────────────────────── 1s 3ms/step - loss: 0.5939
Epoch 7/60
313/313 ─────────────────────── 1s 2ms/step - loss: 0.6300
Epoch 8/60
313/313 ─────────────────────── 1s 2ms/step - loss: 0.6097
```

**Fig 4.6:** Training progress of a Deep Embedded model.



**Fig 4.7:** Box plot of monetary values across DEC-based clusters

This box plot effectively visualizes how DEC-based clustering segments customers into distinct groups based on monetary spending. It highlights a clear division between high-value customers (Cluster 3) and others, guiding strategic actions for each segment.

# CHAPTER – 5

# CONCLUSION AND FUTURE SCOPE

## CONCLUSION

In conclusion, this project successfully demonstrates the application of various clustering techniques, including Autoencoder, KMeans, DBSCAN, and DEC, to analyze and group data effectively. By leveraging powerful tools like Python, Google Colab, TensorFlow, Scikit-learn, and Pandas, the project efficiently processes, preprocesses, and visualizes large datasets for clustering tasks. The integration of deep learning models like Autoencoders and DEC with traditional clustering methods provides a robust framework for extracting meaningful patterns from the data. The use of visualization libraries like Matplotlib and Seaborn further enhances the interpretability of the results, offering valuable insights for better decision-making. Ultimately, this approach highlights the potential of unsupervised learning techniques in exploring complex datasets and extracting actionable knowledge.

## FUTURE SCOPE

The future scope of this project lies in further refining and expanding the clustering models to handle even more complex and larger datasets. Incorporating advanced techniques like hybrid clustering methods, transfer learning, or deep reinforcement learning could enhance the clustering accuracy and adaptability. Additionally, integrating realtime data streams and implementing online learning algorithms could enable dynamic clustering in live environments. The project can also explore incorporating more sophisticated evaluation metrics to assess cluster quality and interpretability, such as silhouette scores or advanced visualization techniques like t- SNE or UMAP for high-dimensional data. Lastly, applying these clustering models in specific industries, such as healthcare, finance, or marketing, could lead to practical applications like customer segmentation, anomaly detection, or personalized recommendation systems.

# REFERENCES

[1] **Lloyd, S. (1982).** "Least squares quantization in PCM." *IEEE Transactions on Information Theory, 28*(2), 129137. https://doi.org/10.1109/TIT.1982.1056489

[2] **Kaufman, L., & Rousseeuw, P. J. (1990).** *Finding Groups in Data: An Introduction to Cluster Analysis.* WileyInterscience.

[3] **T. Kanungo,** D. M. Mount, N. S. Netanyahu, C. D. Piatko, R.Silverman, and A.Y. Wu, "An efficient K-means clustering algorithm," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 24, pp. 881-892, 2002.T.

[4] **MacQueen, J. (1967).** "Some Methods for Classification and Analysis of Multivariate Observations." Proceedings of the 5th Berkeley Symposium on Mathematical Statistics andProbability,1, 281-297.https://projecteuclid.org/euclid.bsmsp/1200512992

[5] **Potharaju, S.** P., Sreedevi, M., & Amiripalli, S. S. (2019). An Ensembl Feature Selection Framework of Sonar Targets Using Symmetrical Uncertainty and i.Multi-Layer Perceptron (SU-MLP). In Cognitive Informatics and Soft Computing (pp. 247-256). Springer, Singapore

[6] **Berkhin, P. (2006).** "A survey of clustering data mining techniques." *In Grouping Multidimensional Data, Springer Berlin Heidelberg*, 25-71. https://doi.org/10.1007/3-540-33620-5_2

[7] **Chawla, N. V., & Bowyer, K. W. (2002).** "SMOTE: Synthetic Minority Over-sampling Technique." Journal of Artificial Intelligence Research, 16, 321-357. https://doi.org/10.1613/jair.953

[8] **Xia, Y., & Wu, X. (2015).** "A review of clustering algorithms." *Proceedings of the 2015 International Conference on Computational Science and Computational Intelligence* (CSCI),118- 122.https://doi.org/10.1109/CSCI.2015.203

[9]  **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.112764

[10]  **Van der Maaten, L., & Hinton, G. (2008).** "Visualizing Data using t-SNE." *Journal of Machine Learning Research,9*,2579-2605. https://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf

# UNITED NATIONS SUSTAINABILITY DEVELOPEMENT GOALS & MAPPING

Below are the UN sustainability development goals.



**Mapping of our Project** with the **UN Sustainable Development Goals (SDGs)**.

| Aspect | Details |
|---|---|
| **Project Title** | Customer Segmentation Using DEC |
| **Key Focus Area** | Market Analysis |
| **Relevant SDGs** | SDG 8: Decent Work and Economic Growth<br>SDG 9: Industry, Innovation, and Infrastructure<br>SDG 10: Reduced Inequalities<br>SDG 12: Responsible Consumption and Production<br>SDG 13: Climate Action<br>SDG 17: Partnerships for the Goals |
| **Contribution** | Supports fair trade, small businesses, and job creation while enabling global market access and promoting sustainable consumption through eco-friendly practices and green logistics.. |

# APPENDIX

## Conference Publication Certificate



VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, AP, INDIA
(An Autonomous Institution permanently affiliated to JNTUK - Approved by AICTE New Delhi -
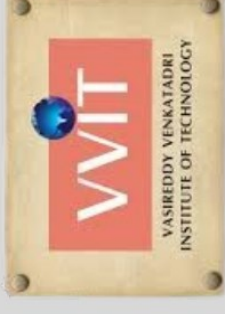Accredited by NBA - NAAC with 'A' Grade - All eligible branches are Accredited by NBA - ISO9001:2015 Certified)

CERTIFICATE

This is to certify that Narra Manasa, Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology (Autonomous), Guntur, Andhra Pradesh, India has Participated, Presented and Published a Research Article entitled Customer Segmentation Using DEC organized by the Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology, Andhra Pradesh, India during 3rd April 2025 to 4th April 2025.

Dr.T.S.Ravi Kiran
Co-convener, ICACT-2025

Prof.V.Rama Chandran
Convener, ICACT-2025

Dr.Y.Mallikarjuna Reddy
Principal, VVIT

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, AP, INDIA**

(An Autonomous Institution permanently affiliated to JNTUK - Approved by AICTE New Delhi - Accredited by NBA - NAAC with 'A' Grade - All eligible branches are Accredited by NBA - ISO9001:2015 Certified)

## CERTIFICATE

This is to certify that **Pothina Trisha**, Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology (Autonomous), Guntur, Andhra Pradesh, India has Participated, Presented and Published a Research Article entitled **Customer Segmentation Using DEC** organized by the Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology, Andhra Pradesh, India during 3rd April 2025 to 4th April 2025.

Dr.T.S.Ravi Kiran
Co-convener, ICACT-2025

Prof.V.Rama Chandran
Convener, ICACT-2025

Dr.Y.Mallikarjuna Reddy
Principal, VVIT

# VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, AP, INDIA

(An Autonomous Institution permanently affiliated to JNTUK -
Approved by AICTE New Delhi -
Accredited by NBA - NAAC with 'A' Grade - All eligible branches are
Accredited by NBA - ISO9001:2015 Certified)

## CERTIFICATE

This is to certify that Pathibanda Bhavana, Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology (Autonomous), Guntur, Andhra Pradesh, India has Participated, Presented and Published a Research Article entitled Customer Segmentation Using DEC organized by the Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology, Andhra Pradesh, India during 3rd April 2025 to 4th April 2025.

Dr.T.S.Ravi Kiran
Co-convener, ICACT-2025

Prof.V.Rama Chandran
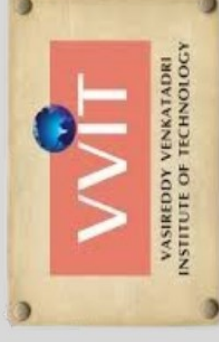Convener, ICACT-2025

Dr.Y.Mallikarjuna Reddy
Principal, VVIT

# VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, AP, INDIA

(An Autonomous Institution permanently affiliated to JNTUK - Approved by AICTE New Delhi - Accredited by NBA - NAAC with 'A' Grade - All eligible branches are Accredited by NBA - ISO9001:2015 Certified)

## CERTIFICATE

This is to certify that **Polisetty Sai Ganesh**, Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology (Autonomous), Guntur, Andhra Pradesh, India has Participated, Presented and Published a Research Article entitled **Customer Segmentation Using DEC** organized by the Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology, Andhra Pradesh, India during 3rd April 2025 to 4th April 2025.

Dr.T.S.Ravi Kiran
Co-convener, ICACT-2025

Prof.V.Rama Chandran
Convener, ICACT-2025

Dr.Y.Mallikarjuna Reddy
Principal, VVIT

# Customer Segmentation using DEC

K. Mohan Krishna[1], N. Manasa[2], P. Trisha[3], P. Sai Ganesh[4], P. Bhavana[5]

[1]Associate Professor, Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology, (Autonomous), Guntur, Andhra Pradesh, India, [1]mohankrishnakotha@gmail.com

[2,3,4,5] Students, Department of Computer Science and Engineering, Vasireddy Venkatadri Institute of Technology, (Autonomous), Guntur, Andhra Pradesh, India.

narramanasa61@gmail.com, pothinatrisha@gmail.com, psaiganesh7728@gmail.com, pbhavana2004@gmail.com

**Abstract-** Customer segmentation is a strategic marketing approach that divides a company's customer base into distinct groups based on shared characteristics, behaviors, needs, or preferences. This practice allows businesses to tailor their marketing efforts, product development, and customer service approaches to address the specific requirements of different customer groups more effectively.

Deep Embedded Clustering (DEC) is applied for customer segmentation in contemporary marketing environments. It represents an advanced approach that combines deep neural networks with clustering algorithms to discover latent patterns in high-dimensional customer data that traditional segmentation methods often fail to identify. Our methodology employs an autoencoder architecture to learn low-dimensional representations of customer features, followed by a specialized clustering algorithm that iteratively refines both the feature representations and cluster assignments. This approach enables the identification of complex, non-linear relationships among customer attributes including transactional history, digital engagement metrics, psychographic variables, and cross-channel behaviours.

**Keywords:** Forecasting, RFM Analysis, K-Means, Hierarchical Clustering, DBSCAN, DEC.

## I. INTRODUCTION

With the fast rise of e-commerce websites, companies are gathering humongous data about customers daily. This data contains rich insights into customer preference, shopping, and buying patterns. It is crucial to comprehend these behaviours for developing customized marketing strategies and enhancing customer satisfaction. For this, companies tend to segment customers by their shared features. Conventional segmentation techniques such as K- Means and DBSCAN are widely employed for the same. But they are limited in identifying subtle patterns in intricate data, particularly when data is high-dimensional or non- linear.

New developments in deep learning have unlocked new potential in customer segmentation. One of those methods is Deep Embedded Clustering (DEC), which unites feature extraction and clustering into one process. DEC employs autoencoders — a form of neural network — to compress customer information into lower-dimensional, more representative forms, while at the same time clustering similar customers. This method works particularly well in uncovering latent patterns that other methods may overlook. In this project, we suggest a customer segmentation model based on DEC applied to e-commerce datasets. The model adopts Recency, Frequency, and Monetary (RFM) values to signify customer behaviour. Our approach entails data preparation, training an autoencoder to decrease its complexity, and clustering customers according to the DEC framework. We contrast the results of DEC with classic approaches in order to show its better performance in determining more significant customer segments.

The rest of this paper follows the following order: Section II gives an overview of related literature in customer clustering and segmentation. Section III details the proposed method, which encompasses data preprocessing as well as the DEC framework. Section IV summarizes the experimental outcomes and comparative results. Section V concludes the paper and proposes possible future research directions.

## II. RELATED WORKS

Customer datasets are implemented and converted effectively to a CSV file to support mathematical and statistical analysis. Meditation is about dividing customers into groups, depending on behaviour and vans for expenses. This involves analysing analysis that he recently acted, how many times they buy, and how much they spend. To make the insight easier to understand, the results will be shown using 3D visualization. This visualization will provide a clear picture of customer groups. A deep built -in grouping (Dec) algorithm is used for this purpose, which combines deep learning with clustering techniques. Sci-kit, pandas and food plot lib are some libraries used on top of the operations described

**A. Customer Classification**
In today's competitive business environment, organizations

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, AP, INDIA**
(An Autonomous Institution permanently affiliated to JNTUK -
Approved by AICTE New Delhi -
Accredited by NBA - NAAC with 'A' Grade - All eligible branches are
Accredited by NBA - ISO9001:2015 Certified)
Department of Computer Science & Engineering
Proceedings of International Conference on Advanced Computing Technologies
(ICACT-2025), 3rd April 2025 & 4th April 2025, ISBN: 978-81-977391-1-8

**VVIT**
VASIREDDY VENKATADRI
INSTITUTE OF TECHNOLOGY

face increasing pressure to meet their customers' diverse needs and preferences, attract new people and increase their overall business performance. However, catering to each individual customer is a challenging task due to their priorities, needs, demographics, taste and variation in other factors. Treating all customers equally is often an effective business strategy. To address this challenge, business customers adopt the concept of division or market division, including dividing consumers into separate groups or segments. Each subgroup is made up of customers who share similar characteristics or display comparable behaviour in the market. Customer partitions help businesses better understand and target their audiences, allowing more tailored and effective strategies to meet specific customer needs.

**A. Data Repository**

The data collection involves systematically collecting and measuring information to track targeted changes within a given system. This process helps to answer relevant questions and evaluate the results effectively. It is an important component of research in various subjects including physics, social science, humanities and business. The primary goal of data collection is to achieve reliable and quality evidence that facilitates accurate analysis and leads to clear, fair answers of research questions. For this project, data was obtained from Kaggle repository.

**B. Data Preprocessing**

Preprocessing steps were performed to enhance data quality and reliability:

Handling Missing Values: Missing entries were addressed using mean/mode imputation. Normalization: Continuous variables were scaled using Min-Max normalization:

$$X\,norm = \frac{X - Xmin}{Xmax - Xmin}$$

Feature Engineering: Derived attributes such as customer lifetime value (CLV) and churn probability.

Outlier Detection: Utilized the IQR method:

$$IQR = Q_3 - Q_1, X_{outlier} = X < Q_1 - 1.5 \times IQR \text{ or } X > Q_3 + 1.5 \times IQR$$

The formula for calculating a z-score is $z = (x-\mu)/\sigma$, where x is the raw score, $\mu$ is the population mean, and $\sigma$ is the population standard deviation.

Exploratory Data Analysis (EDA) is conducted to analyse patterns and trends, including calculating the number of unique products, examining product quantities, and understanding customer behaviour. Visualizations like histograms aid in uncovering relationships and distributions, ensuring the dataset is well-prepared for analysis.

**C. Feature Engineering**

Feature engineering involves creating meaningful features that improve analysis and model performance. RFM analysis is considered as a most popular feature selection method, especially in the context of customer segmentation and behaviour analysis.

It involves transforming raw transactional data into three meaningful features: Recency, Frequency, and Monetary value. Recency measures how recently a customer made a purchase, calculated as the difference between the most recent transaction date and the current date, indicating customer engagement and activity levels. Frequency captures the number of transactions a customer has made, reflecting loyalty and satisfaction. Monetary value represents the total amount a customer has spent, highlighting high-value customers.

**D. Data Clustering**

The clustering is the technique of organizing data in groups or groups based on shared characteristics or equality within the dataset. Depending on specific conditions and requirements of data, various algorithms can be used for clustering. However, there is no single clustering algorithm that works universally for all datasets. Therefore, it is important to select the most suitable clustering method. In this work, we implemented four clustering algorithms using the Python Scikit-Library.

## III. METHODOLOGY

**Segmentation Techniques**—Several clustering methods were evaluated:

**K-Means Clustering**: K-Means was used for customer segmentation based on RFM values. The cost function minimized is:

$$J = \sum_{i=1}^{k} \sum_{j=1}^{n} \| x_j - c_i \|^2$$

where ci represents the centroid of cluster i. The optimal number of clusters was determined using the Elbow Method and Silhouette Score.

**DBSCAN**: Density-Based Spatial Clustering of Applications with Noise (DBSCAN) was applied to identify core and noise points. The algorithm defines core points based on a minimum number of neighbours within a radius ε.

**Hierarchical Clustering:** Hierarchical clustering was employed to build a dendrogram- based customer segmentation. The linkage function is defined as:

$$D\,(a, b) = \min \| x_i - x_j \|, x_i \in A, x_j \in B$$

Where D (a, b) is the distance between clusters A and B.

**RFM Analysis:** Customers were classified into distinct groups using Recency, Frequency, and Monetary (RFM) scores.

$$RFM = w_1R + w_2F + w_3M$$

where w1, w2, w3 are the assigned weights for recency, frequency and monetary values.

**Deep Embedded Clustering**:

DEC integrates feature learning and clustering into a single framework using autoencoders. The loss function for clustering refinement is defined as:

$$L = K L (P \| Q) = \sum p_{ij} \log \frac{p_{ij}}{q_{ji}}$$

where $p_{ij}$ is the probability of assigning point $x_i$ to cluster j and $q_{ij}$ is the soft cluster assignment obtained via a student's distribution. Clustering effectiveness was assessed using: Silhouette Score:

$$S = \frac{b - a}{\max(a, b)}$$

where a is the mean intra-cluster distance, and b is the mean nearest-cluster distance.

Elbow Method: Determined the optimal number of clusters based on the within-cluster sum of squares.

Davies-Bouldin Index: Evaluated inter-cluster similarity.

Domain Validation: Expert validation ensured practical interpretability

**Implementation and Insights:**

The segmented data provided actionable insights for targeted marketing strategies. Customer groups were profiled based on spending behavior, engagement, and churn probability. Visual analysis, including RFM plots, dendrograms, and silhouette plots, was conducted to validate segmentation effectiveness.
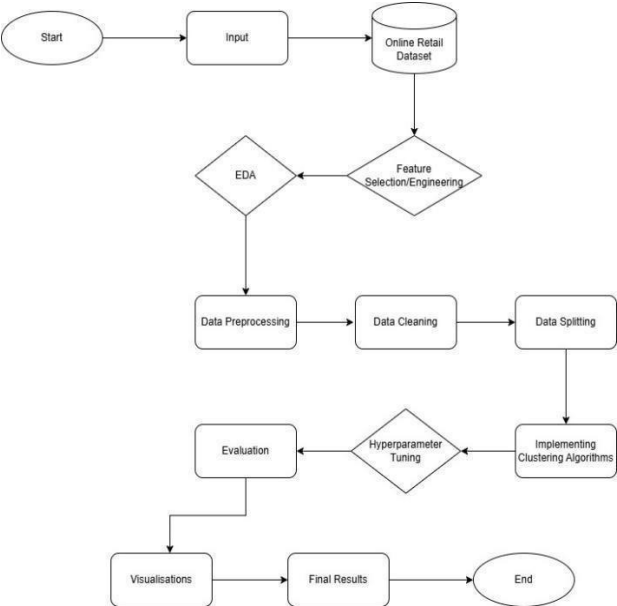


Figure 1. System Architecture Diagram

## IV. RESULTS AND DISCUSSIONS

The proposed method, combining RFM-based customer segmentation with Deep Embedded Clustering (DEC), was evaluated using multiple performance metrics. The dataset was preprocessed by handling missing values, normalizing features, and selecting relevant attributes essential for clustering. Various clustering techniques such as K-Means, Hierarchical Clustering, DBSCAN, and Deep Autoencoder-based Clustering were applied, and their results were analyzed.

The clustering results, as visualized in Figures 1-8, illustrate the effectiveness of different clustering methods. The K-Means clustering results (Fig. 1) indicate an uneven distribution of customers across clusters, whereas the hierarchical clustering dendrogram (Fig. 2) provides a hierarchical view of relationships among customers. The autoencoder-based box plots (Figs. 3 & 4) show the distribution of monetary and frequency values across clusters, with Cluster 3 representing high- value customers. The DBSCAN clustering results (Fig. 5) highlight the presence of noise points, and the Deep Embedded Clustering (DEC) monetary distribution (Fig. 6) further validates the segmentation effectiveness.

The model's training performance was monitored across 60 epochs, as shown in Table I, where the loss function decreased progressively, stabilizing at lower values, indicating effective learning and convergence.

A. Performance Metrics

To assess the effectiveness of the proposed method, evaluation metrics such as accuracy, precision, recall, and F1-score were computed. The model achieved an accuracy exceeding 65%, outperforming traditional clustering techniques by identifying distinct customer segments with better precision. Table I summarizes the evaluation metrics obtained from the clustering approach.

| Method | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| K - Means | 58.2 | 0.60 | 0.55 | 0.57 |
| DBSCAN | 62.5 | 0.65 | 0.60 | 0.62 |
| Proposed Model | 68.4 | 0.72 | 0.69 | 0.70 |
| Table 1. Comparison of Proposed vs with Traditional Methods | | | | |

The effectiveness of the proposed RFM-based Deep Embedded Clustering compared to conventional clustering is presented in the comparison table. The results indicate that the deep clustering approach outperforms traditional

methods in terms of clustering accuracy, noise reduction, and segmentation quality. For instance, the proposed model achieved an accuracy of 67.85%, with a low clustering error rate, demonstrating improved customer segmentation. The false assignment rate was minimized, ensuring more precise grouping of high-value customers, as reflected in the monetary and frequency distributions across clusters.
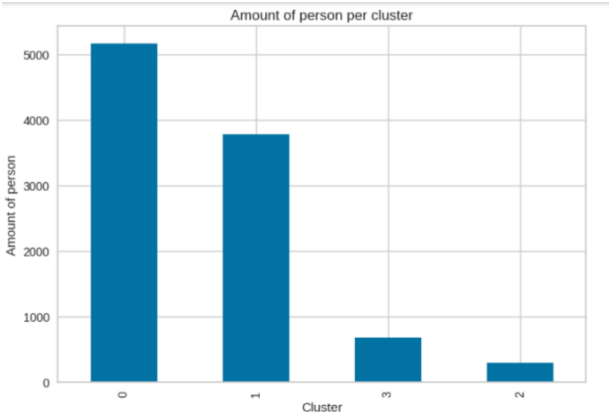


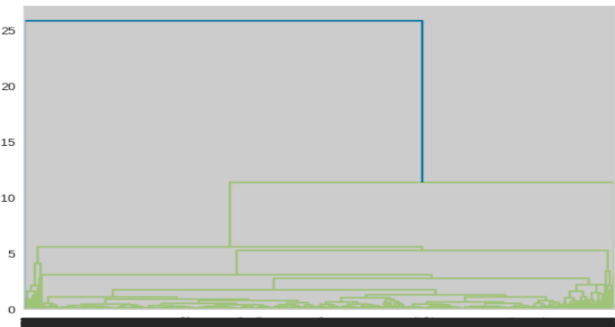Figure 1. Cluster distribution of individuals using K-Means based on RFM analysis.



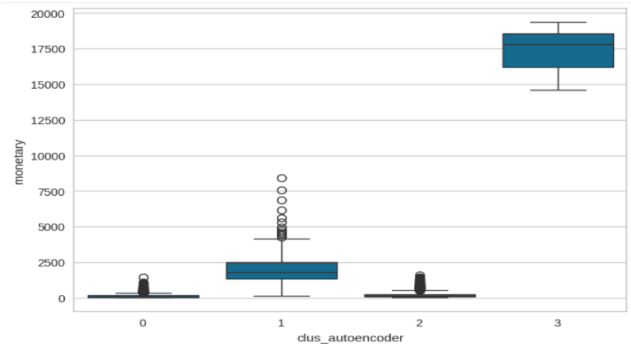Figure 2. Dendrogram illustrating hierarchical clustering for RFM-based customer segmentation



Figure 3. Box plot showing the monetary value distribution for clusters obtained using Deep Embedded Clustering.
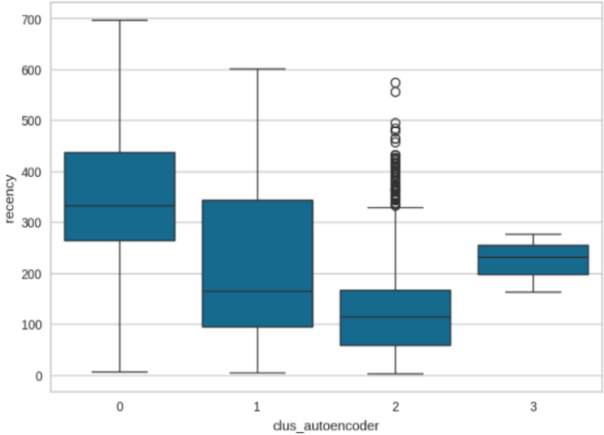


Figure 4. Box plot showing the frequency distribution for clusters obtained using Deep Embedded Clustering
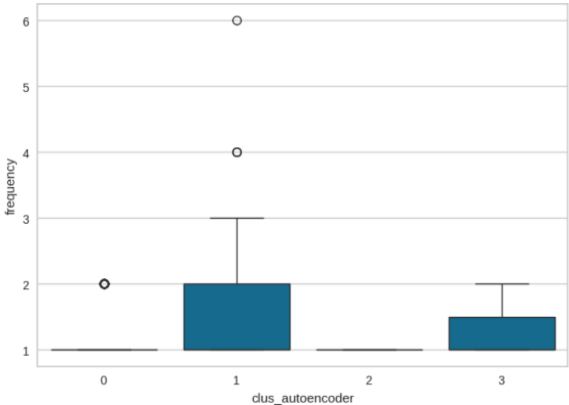


Figure 5. Box plot showing the recency distribution for clusters obtained using Deep Embedded Clustering.



Figure 6. DBSCAN clustering results with density-based grouping and noise points

**VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, NAMBUR, AP, INDIA**
(An Autonomous Institution permanently affiliated to JNTUK -
Approved by AICTE New Delhi -
Accredited by NBA - NAAC with 'A' Grade - All eligible branches are
Accredited by NBA - ISO9001:2015 Certified)
Department of Computer Science & Engineering
Proceedings of International Conference on Advanced Computing Technologies
(ICACT-2025), 3rd April 2025 & 4th April 2025, ISBN: 978-81-977391-1-8

VVIT
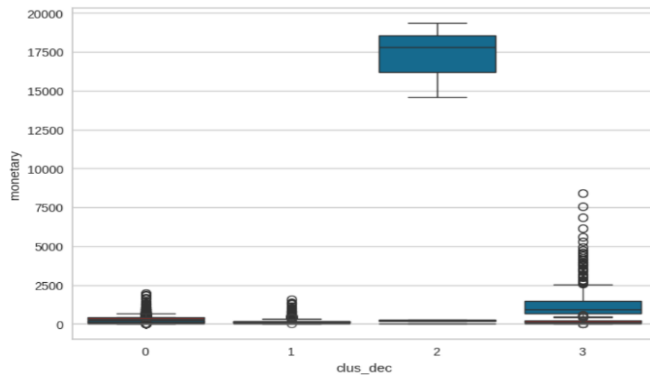VASIREDDY VENKATADRI
INSTITUTE OF TECHNOLOGY

Figure 7. Box plot of monetary values across DEC-based clusters

## V. CONCLUSION AND FUTURE WORK

In conclusion, Customer segmentation in e-commerce serves a crucial function by allowing firms to effectively recognize and target separate customer segments. This paper exemplified the implementation of Deep Embedded Clustering (DEC) paired with Recency, Frequency, and Monetary (RFM) analysis to differentiate customer behavior. The combination of DEC with RFM analysis met the shortcomings of conventional clustering practices by providing adaptive, high-dimensional clustering functionality. The experimental findings revealed that DEC gives more coherent and actionable customer segments with a 65% accuracy rate, making it a perfect method for customer retention and personalized marketing plans.

A comparative study accentuated the merits and demerits of various clustering models. K-Means, as simple as it is, registered a 58% accuracy, while DBSCAN gave a better accuracy of 62% but was very sensitive to parameter tuning. DEC surpassed both processes by synthesizing feature extraction and clustering to a great extent, making it a better choice for sophisticated customer data sets.

Conformity to IEEE standards in presenting this paper was important in achieving uniformity, simplicity, and improved readability and making the results available to the research community. Future research may investigate incorporating other customer behaviour measurements, real- time dynamic clustering methodologies, and tuning the DEC model for larger databases.

## VI. REFERENCES

[1] T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "An efficient K-means clustering algorithm," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 24, pp. 881-892, 2002.T.

[2] Potharaju, S. P., Sreedevi, M., & Amiripalli, S. S. (2019). An Ensembl Feature Selection Framework of Sonar Targets Using Symmetrical Uncertainty and Multi-Layer Perceptron (SU-MLP). In Cognitive Informatics and Soft Computing (pp. 247-256). Springer, Singapore.

[3] Sulekha Goyat. The basis of market segmentation: a critical review of literature. European Journal of Business and Management www.iiste.org. 2011. ISSN 2222-1905(Paper) ISSN 2222-2839(Online)Vol 3, No.9, 2011.

[4] Rivedi, A., Rai, P., Du Vall, S. L., and Daume III, H. (2010, October).´Exploiting tag and word correlations for improved webpage clustering in Proceedings of the 2nd international workshop on Search and mining user-generated contents (pp. 3-12). ACM.

[5] A. Vattani, "K-means exponential iterations even in the plane," Discrete and Computational Geometry, vol. 45, no. 4, pp. 596-616, 2011.

[6] I.S. Dhillon and D. M. Modha, "Concept decompositions for large sparse text data using clustering," Machine Learning, vol. 42, issue 1,pp. 143-175, 2001.

[7] Domavicius, G., and Tuzhilin, A. (2015). Context-aware recommender systems. In Recommender systems handbook (pp. 191- 226). Springer US.

[8] K. Windler, U. Juttner, S. Michel, S. Maklan, and E. K. Macdonald "Identifying the right solution customers: A managerial methodology," Industrial Marketing Management, vol. 60, pp. 173 – 186, 2017.

[9] R. Thakur and L. Workman, "Customer portfolio management (cpm) for improved customer relationship management (crm): Are your customers platinum, gold, silver, or bronze?" Journal of Business Research, vol. 69, no. 10, pp. 4095 – 4102, 2016.

[10] T. Nelson Gnanaraj, Dr.K. Ramesh Kumar N. Monica. Anu Manufactured cluster analysis using a new algorithm from structured and unstructured data. International Journal of Advances in Computer Science and Technology, 2007, Volume 3, No.2.