# Advanced Market Segmentation Using Deep Clustering

## Phase 1: Problem Definition and Data Understanding

### 1.1 Project Overview

The primary objective of this project is to implement advanced market segmentation using deep clustering techniques. Market segmentation is essential for identifying distinct groups of customers based on their behavior, preferences, or demographics. Traditional clustering methods struggle with high-dimensional data, leading to suboptimal segmentation. To address this, we leverage deep learning, specifically auto encoders, to extract meaningful latent features and enhance the clustering process.

The goal is to divide customers into distinct segments using unsupervised learning, which can then be targeted with tailored marketing strategies, personalized recommendations, or improved customer service. The project aims to provide businesses with a deep understanding of customer groups, leading to more informed decision-making.

### 1.2 Objective of the Project

- **Objective**: The objective of this project is to implement a **clustering** model that can identify distinct customer segments based on their behaviors, demographics, and engagement with products or services.

  **Clustering** is the core objective of this project because it involves grouping similar data points (customers) together without prior labels or classifications. This is an unsupervised learning task where the goal is to automatically discover patterns and relationships within the data.

- **Target Users**: This project is primarily aimed at businesses and marketers who want to gain insights into customer behaviors and preferences. It is also valuable for data scientists and machine learning practitioners interested in applying deep learning techniques to clustering problems.
- **Potential Applications**:
    - **Customer Segmentation**: Businesses can use the model to group customers with similar behaviors and preferences, enabling targeted marketing campaigns, personalized recommendations, and improved customer service.
    - **Product Development**: Identifying customer segments can inform product design and feature prioritization by focusing on the needs and preferences of different groups.

- o **Customer Support**: Segments can help support teams provide tailored assistance, addressing common issues that arise within each customer group.

## 1.3 Dataset Overview and Data Requirements

To achieve the goal of customer segmentation, the dataset needs to include features related to customer behaviors, demographics, and engagement with products or services. The dataset format must support both categorical and continuous data types to enable comprehensive analysis.

- **Features**:
  - o **Demographics**: Information such as age, gender, income, occupation, and geographic location.
  - o **Behavioral Features**: Data related to customer purchases, frequency of transactions, types of products bought, and amount spent.
  - o **Engagement Features**: Interaction data including clicks, visits to websites, responses to marketing campaigns, and social media activity.
  - o **Additional Features**: Any other customer data that could influence purchasing decisions, such as time spent on the website or customer feedback scores.
- **Labels**: Since this is an unsupervised learning task, there are no explicit labels in the dataset. The objective is to automatically group the data based on the relationships between features, without predefined categories.
- **Dataset Format**:
  - o The data should be in tabular format (e.g., CSV, Excel, or SQL database).
  - o Each row represents an individual customer, with columns for various customer attributes and behaviors.
  - o The dataset may also include timestamps or categorical data (e.g., product categories, customer segments) that need to be appropriately encoded for machine learning tasks.

## 1.4 Data Sources

The data required for this project can be sourced from various locations, both public and proprietary. The following are possible sources for customer data:

- **Public Datasets**:
  - o **UCI Machine Learning Repository**: The repository includes datasets for customer behavior and market segmentation that can be leveraged to build initial models.

- o **Kaggle Datasets**: Kaggle offers several publicly available datasets related to customer segmentation, such as customer behavior data from online retail stores or financial institutions.
  - o **Google Dataset Search**: A comprehensive search tool that indexes public datasets on various domains, including market segmentation.
- **Web Scraping**:
  - o **E-commerce websites**: Data can be scraped from e-commerce platforms like Amazon, eBay, or local online retailers to gather information about customer purchases, product preferences, and behaviors.
  - o **Social Media**: Social media platforms such as Twitter or Instagram can provide engagement data, where scraping can be done to analyze customer interactions with brand-related content.
- **Proprietary Data**:
  - o **Company CRM Systems**: Businesses often collect detailed customer data through their customer relationship management (CRM) systems. This can include purchase histories, demographic details, and customer feedback.
  - o **Sales and Marketing Data**: Customer purchase and interaction data from internal company sales systems, loyalty programs, or marketing campaigns can be a rich source of insights for segmentation.

## 1.5 Initial Data Exploration

Once the dataset has been sourced, an initial data exploration phase will be conducted to understand the quality and structure of the data. The tasks involved in this phase include:

- **Missing Data**: Identifying columns with missing data and applying imputation strategies, such as mean imputation for numerical data or mode imputation for categorical data.
- **Outliers**: Outlier detection and treatment to ensure that extreme values do not negatively impact the performance of the clustering algorithms.
- **Data Distribution**: Analyzing the distribution of key features to determine if any transformations (e.g., normalization or scaling) are required to ensure that the data is ready for deep learning techniques.
- **Correlation Analysis**: Identifying correlations between features to help understand the relationships in the dataset and to assist in feature selection or reduction.
- **Exploratory Visualizations**: Using histograms, scatter plots, and pair plots to visualize the data and identify any patterns or trends that can inform the next steps in model development.

## 1.6 Preprocessing Objectives

The goal of preprocessing is to transform raw data into a format that can be effectively used for model development. This includes:

- **Feature Scaling**: Applying scaling techniques such as Min-Max scaling or Z-score normalization to ensure that all numerical features have similar scales.
- **Categorical Encoding**: Converting categorical variables into numerical representations using techniques like one-hot encoding.
- **Feature Selection**: Removing irrelevant or highly correlated features to reduce dimensionality and ensure that only meaningful features are used in the model.
- **Data Transformation**: Applying log transformations, or other techniques, if necessary, to deal with skewed or non-linear features.

**1.7 Conclusion of Phase 1**

Phase 1 has provided a comprehensive understanding of the project's objectives, data requirements, and the sources of data that will be used. The dataset, containing both demographic and behavioral features, will be preprocessed and explored to prepare it for deep learning techniques in the subsequent phases. The goal of market segmentation will be achieved by applying advanced deep clustering methods, which will be the focus of the next phases of the project.