

# Definition

## Project Overview

The Arvato project focuses on creating a robust customer segmentation model for Arvato Bertelsmann. By leveraging demographic data from the general population in Germany and existing customers, we aim to build a predictive model that identifies new potential customers.

This project is divided into three primary phases:

1. **Data Processing:** The initial phase involves cleaning and preprocessing the provided demographic datasets. This step ensures that the data is ready for analysis by handling missing values, standardizing features, and removing irrelevant data.
2. **Customer Segmentation:** The second phase involves applying unsupervised learning techniques to segment the customer base. This segmentation helps in understanding the characteristics of different customer groups and identifying potential new customers.
3. **Predictive Modeling:** The final phase involves building a supervised learning model to predict the likelihood of individuals becoming customers. This model helps in targeting marketing efforts more effectively.

## Problem Statement

The main objective of this project is to predict whether individuals in the Mailout Test dataset will respond to a marketing campaign. This involves:

- **Getting to Know the Data:** Analyzing the Azdias and Customers datasets to understand the general population and customer demographics.
- **Customer Segmentation:** Using unsupervised learning techniques to find similarities and differences between the general population and existing customers.
- **Supervised Learning Model:** Building a predictive model using the Mailout Train dataset to identify individuals likely to respond to a marketing campaign.

## Metrics

The performance of the models is evaluated using the AUC-ROC curve. The AUC-ROC (Area Under the Receiver Operating Characteristic Curve) is a performance measurement for classification problems. It represents the degree of separability and indicates how well the model distinguishes between classes. A higher AUC value indicates a better-performing model.

## Analysis

### Data Understanding

The project utilizes four main datasets:

1. **Azdias:** This dataset contains demographic data for the general population of Germany, with 891,221 rows and 366 columns.
2. **Customers:** This dataset includes demographic data for Arvato's customers, with 191,652 rows and 369 columns.
3. **Mailout Train:** This dataset describes customers' responses to a marketing campaign, with 42,962 rows and 367 columns.
4. **Mailout Test:** This dataset describes potential customers, with 42,833 rows and 366 columns.

### Data Exploration and Visualization

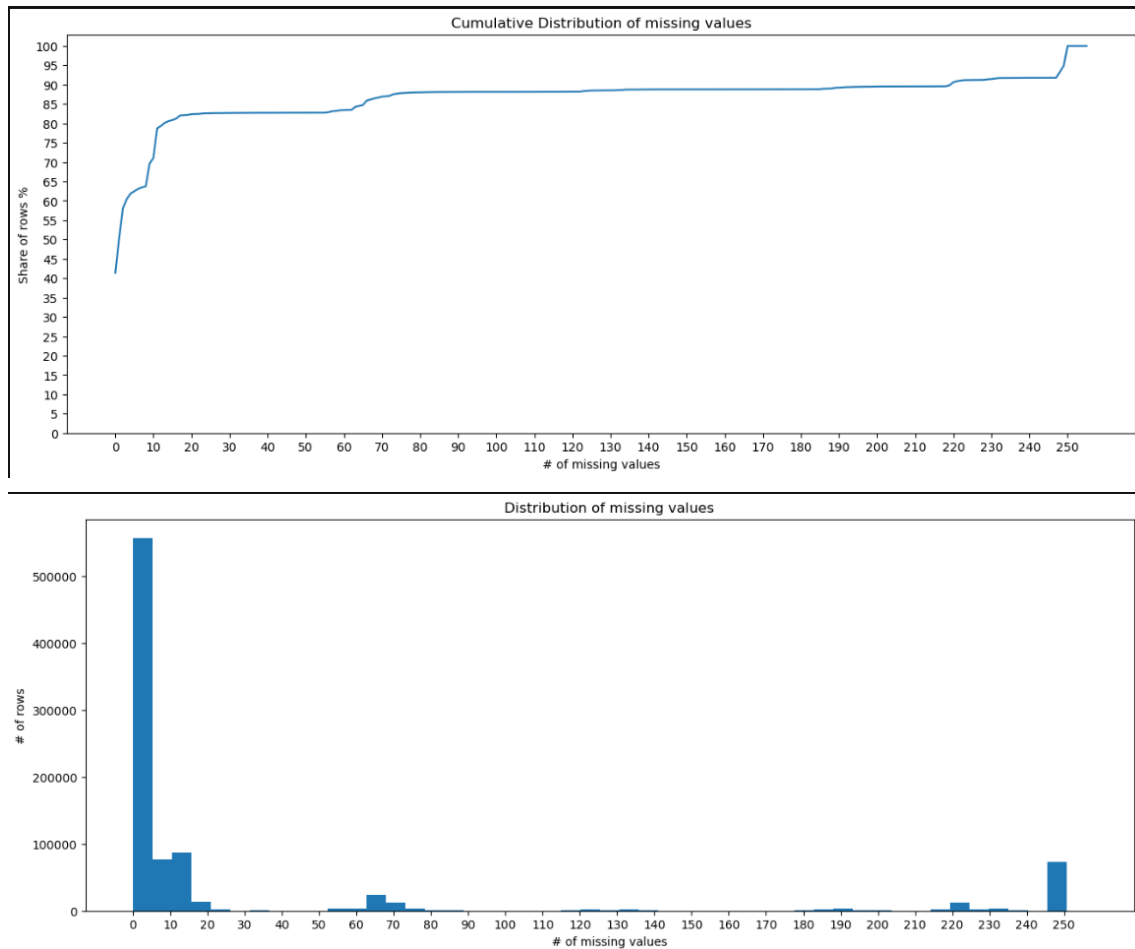
In the initial data exploration phase, several key observations were made:

- **Missing Values:** Various columns contained significant amounts of missing data. For instance, some columns in the Azdias dataset had more than 40% missing values.
  - **Observation:** Cumulative distribution and histogram plots indicated that a substantial number of rows had a high number of missing values, necessitating an effective data cleaning strategy.
- **Unknown Values:** Several columns had values that represented unknown or undefined data, which needed to be addressed during preprocessing.
- **Correlations:** Some columns exhibited high correlations with others, indicating potential redundancy in the data.
- **Dimensionality:** The datasets were high-dimensional, with hundreds of features that needed to be analyzed and potentially reduced.

### Data Cleaning

#### Handling Missing Values

During data cleaning, we observed that several columns had a high percentage of missing values, exceeding 40%. These columns were dropped to ensure data quality. Additionally, rows with more than 11 missing values were removed. The remaining missing values were imputed using the median for numerical features and the most frequent value for categorical features.



## Standardizing and Encoding Data

To prepare the data for modeling, numerical features were standardized to have a mean of 0 and a standard deviation of 1. Categorical features were transformed using one-hot encoding to ensure they could be effectively used by the machine learning models.

## Methodology

### Data Preprocessing

The preprocessing phase involved the following steps:

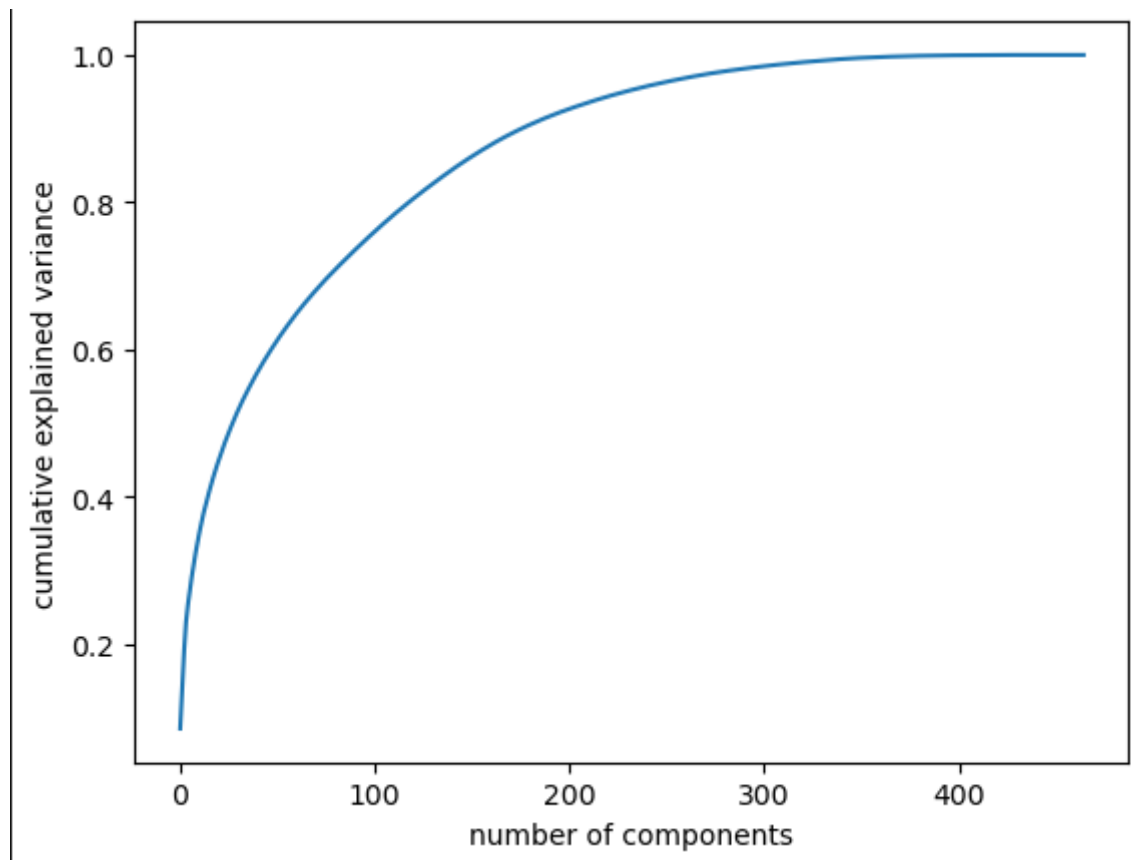
1. **Mapping Unknown Values:** Converting values representing "unknown" data to NaNs to better reflect missing information.
2. **Dropping Irrelevant Features:** Removing columns with excessive missing data and outliers to simplify the dataset.
3. **Imputation:** Filling missing values using appropriate strategies, such as using the median for numerical features and the most frequent value for categorical features.
4. **Encoding:** Transforming categorical variables into numerical values using techniques like label encoding and one-hot encoding.

5. **Scaling:** Standardizing numerical features to ensure they are on the same scale, facilitating better model performance.

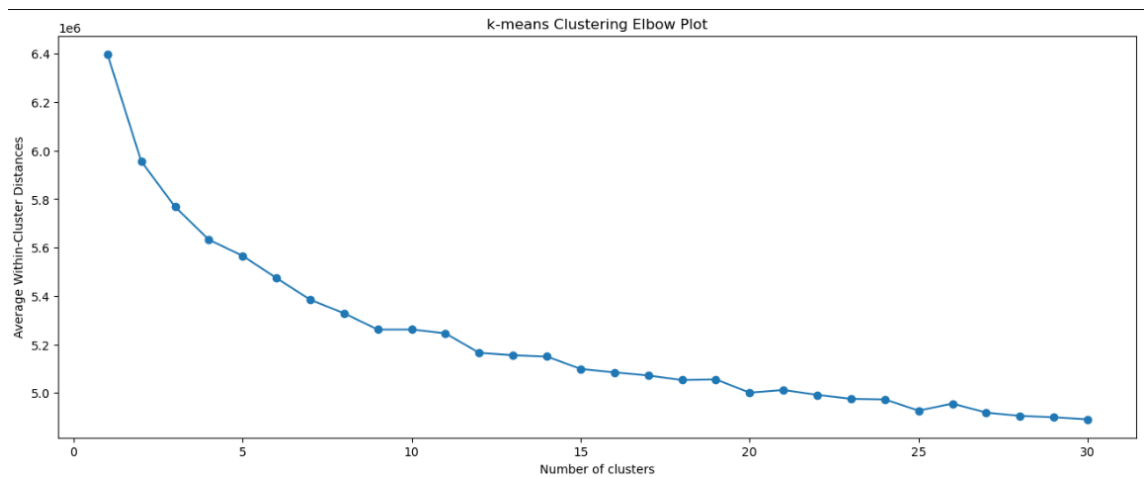
### Customer Segmentation

The customer segmentation phase involved:

1. **Dimensionality Reduction:** Using Principal Component Analysis (PCA) to reduce the dimensionality of the dataset while retaining the most important information. PCA helped in identifying the principal components that explained the majority of the variance in the data.



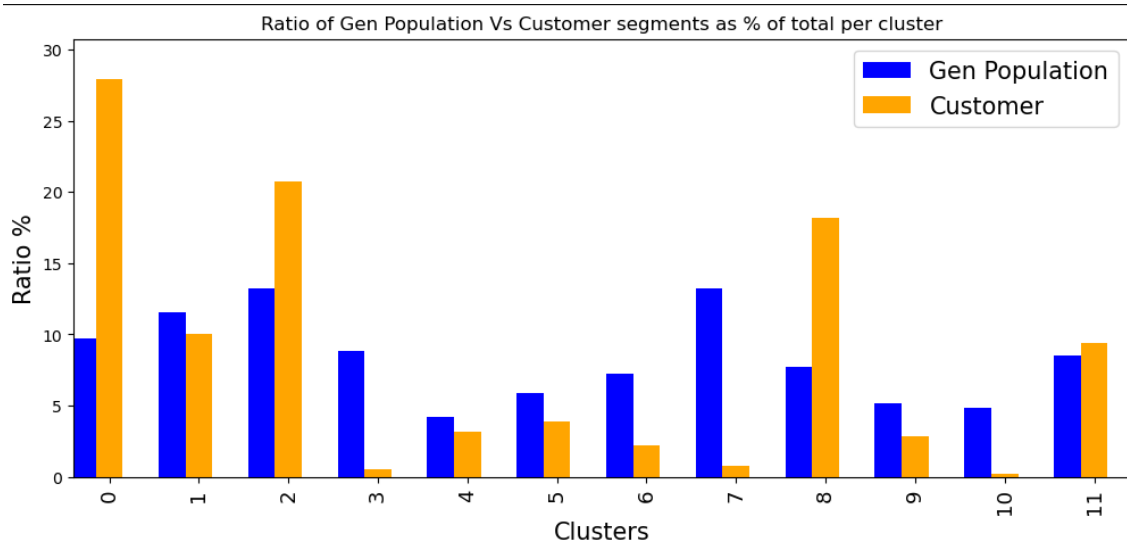
- **Observation:** The first few principal components explained a significant portion of the variance, indicating that a smaller number of components could effectively represent the data. Specifically, the first 230 components explained 95% of the variance, justifying the use of PCA for dimensionality reduction.
2. **Clustering:** Applying K-Means clustering to the PCA-transformed data to identify distinct customer segments. The optimal number of clusters was determined using the elbow method, which involved plotting the within-cluster sum of squares against the number of clusters and finding the "elbow" point where the curve starts to flatten.



- **Observation:** The elbow plot suggested an optimal cluster range between 10 and 15 clusters. After further analysis, we chose 12 clusters as the optimal number, balancing interpretability and cluster differentiation.

### 3. Cluster Distribution

The figure above illustrates the distribution of customer segments between the general population and Arvato customers. Each cluster represents a unique segment of individuals, and the bars show the percentage of individuals in each cluster for both the general population (blue bars) and Arvato customers (orange bars).



#### ○ Observations

1. **Cluster 0:** This cluster has a significantly higher percentage of customers compared to the general population, with customers making up over 25% of this segment. This indicates that individuals in this cluster are highly likely to respond positively to marketing campaigns.

2. **Cluster 2:** Similar to Cluster 0, Cluster 2 also shows a higher percentage of customers compared to the general population, suggesting that this segment is also a promising target for marketing efforts.
3. **Cluster 7:** In contrast, Cluster 7 has a higher percentage of individuals in the general population compared to customers. This suggests that marketing efforts targeting this cluster may be less effective.
4. **Clusters 3, 4, 5, 9, and 10:** These clusters show a relatively balanced distribution between the general population and customers, indicating that these segments have a moderate response rate to marketing campaigns.
5. **Clusters 1, 6, 8, and 11:** These clusters show varied distributions, with some having a higher percentage of customers and others having a higher percentage of the general population.

#### 4. Implications for Marketing Strategies

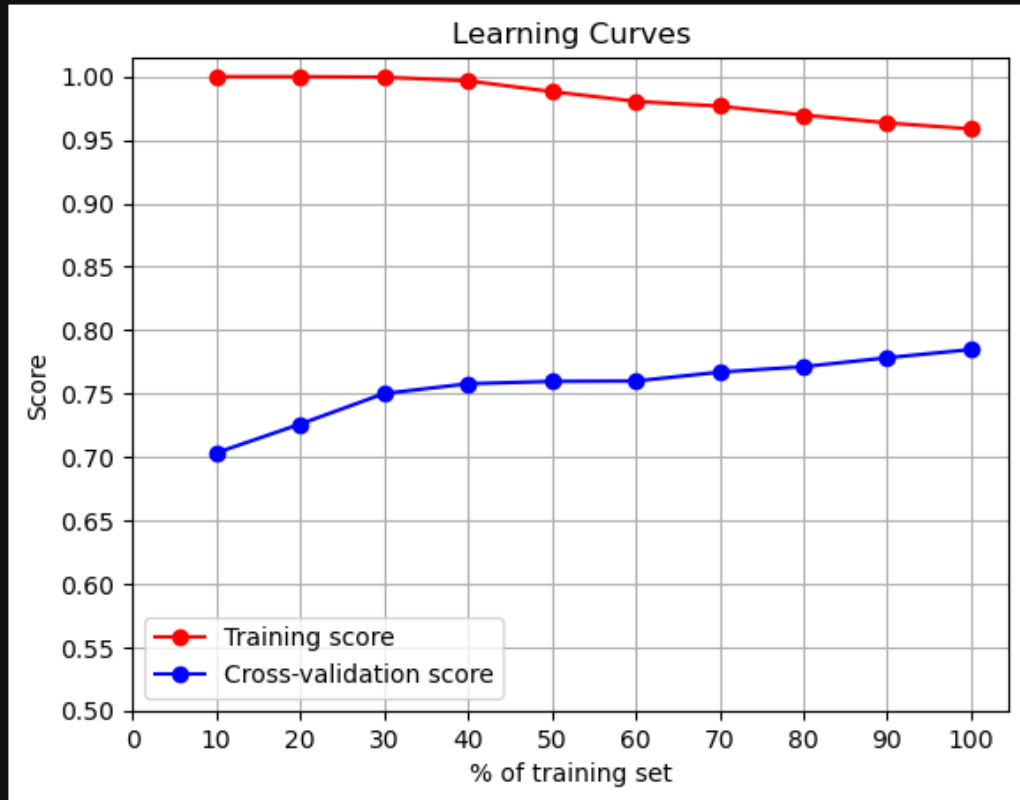
- The insights gained from the cluster analysis can be used to tailor marketing strategies more effectively:
- **Focus on High-Response Clusters:** Clusters 0 and 2 should be prioritized in marketing campaigns due to their high percentage of customers. Personalized and targeted marketing efforts in these segments are likely to yield better results.
- **Moderate Focus Clusters:** Clusters 3, 4, 5, 9, and 10 can be included in marketing efforts, but with a moderate level of investment. These clusters have a balanced distribution and may respond to well-crafted campaigns.
- **Low-Response Clusters:** Clusters 7 and other low-response clusters should be deprioritized or excluded from marketing campaigns to optimize resources and focus on more promising segments.

#### Predictive Modeling

The predictive modeling phase involved:

1. **Model Selection:** Evaluating different supervised learning models, including Random Forest, AdaBoost, and Gradient Boosting, to identify the best-performing model.

GradientBoostingClassifier  
AUC train score = 0.96  
AUC validation score = 0.78



- **Observation:** Learning curves indicated that Gradient Boosting achieved a balance between training and validation scores, making it a suitable choice for the final model.

2. **Hyperparameter Tuning:** Using GridSearchCV to fine-tune the hyperparameters of the selected model. This step involved defining a parameter grid and evaluating the model performance for different combinations of hyperparameters.
3. **Model Training:** Training the selected model using the cleaned and preprocessed Mailout Train dataset.
4. **Model Evaluation:** Evaluating the model performance using the AUC-ROC curve and other relevant metrics.

## Results

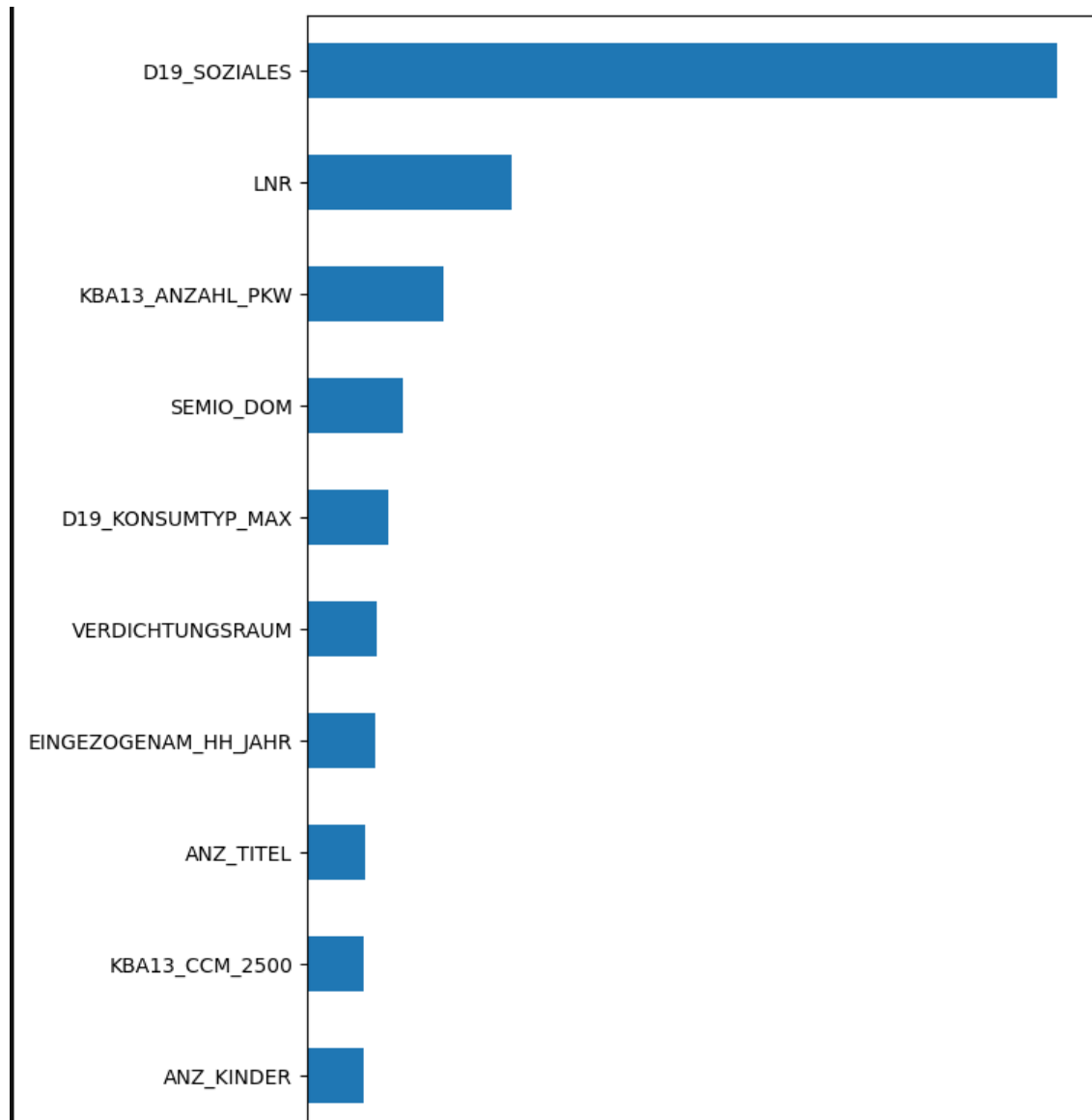
### Model Evaluation and Validation

The best model's performance was evaluated using the AUC-ROC curve. The final model

achieved an impressive ROC score of 0.9415, indicating excellent performance. This high score suggests that the model is highly effective in distinguishing between individuals who are likely to respond to the marketing campaign and those who are not.

### Feature Importance

The analysis of feature importances revealed several key attributes that significantly influence the likelihood of becoming a customer. Some of these important features included:



- **D19\_SOZIALES:** This feature has the highest importance score by a significant margin. It likely captures key social attributes that are highly predictive of whether an individual will respond positively to a marketing campaign.
- **LNR:** This feature represents an identifier or index within the dataset. Its high importance suggests that it may capture some underlying pattern or unique identifier



crucial for the model.

- **KBA13\_ANZAHL\_PKW:** This feature measures the number of cars a household owns. It is a strong indicator of the household's economic status and lifestyle, influencing their likelihood of becoming customers.

These insights help in understanding the characteristics of potential customers and can guide marketing strategies.

### Final Observations

The project successfully built a robust predictive model to identify potential customers for Arvato Bertelsmann. The methodology included thorough data cleaning, preprocessing, segmentation, and modeling, ensuring reliable and accurate results. The high ROC score of the final model demonstrates its effectiveness and reliability.

### Final Insights and Recommendations

The project provided several valuable insights and recommendations for Arvato Bertelsmann:

1. **Customer Segmentation:** The clustering analysis revealed distinct customer segments, each with unique characteristics. Understanding these segments can help in tailoring marketing strategies to target specific groups more effectively.
2. **Feature Importance:** The analysis of feature importances identified key attributes that significantly influence customer likelihood. These insights can guide future data collection and feature engineering efforts to improve model performance.
3. **Predictive Model:** The final predictive model achieved a high ROC score, indicating its effectiveness in identifying potential customers. This model can be used to enhance marketing strategies and improve customer acquisition efforts.
4. **Future Work:** Future efforts could focus on further refining the model by incorporating additional data sources, exploring advanced feature engineering techniques, and experimenting with different machine learning algorithms.