

# Capstone Report

## Customer Segmentation and Optimization of Customer Acquisition with Arvato Financial Solutions

Rodrigo Polverari Pacheco de Toledo

9<sup>th</sup>March2022

Machine Learning Engineer Nanodegree  
School of Artificial Intelligence

# Contents

I. Definition	3
Project Overview	3
Datasets and Inputs	3
Problem Statement	4
Evaluation Metrics	4
II. Analysis	5
Data Exploration	5
Data Pre-processing	5
III. Algorithms, Techniques and Methodology	9
Customer Segmentation	9
1. Dimensionality Reduction	9
2. Clustering	10
Customer Acquisition	12
IV. Results	14
V. Justification	15
VI. Improvements and Future Steps	15

# I. Definition

## Project Overview

Arvato, wholly owned by Bertelsmann, is a services company that actively develops and implements innovative solutions for customers on a global scale. It provides services including customer support, information technology, logistics and finance [1]. As part of its services, Arvato is helping client companies get invaluable insights into client profiling and marketing.

In this project, we will employ the use of machine learning to deal with real-life data provided by Bertelsmann Arvato Analytics. More specifically, Arvato is helping a client mail-order company to better target next probable customers for its products.

To fulfill the goal, we will focus on customer segmentation, that is, characterizing customers segment of population based upon well-defined specific features [2]. Diving customers into groups based on common characteristics will enable our client company to market each group effectively and properly. Here, we will analyze demographics data of customers and the general German population. This will be followed by developing a supervised machine learning model to make predictions on whether a person will be a new customer.

Access to the comprehensive project files is given here:

<https://github.com/ststapy-ml/Arvato-ML>.

## Datasets and Inputs

The project makes use of four data files:

- 1) Udacity\_AZDIAS\_052018.csv: Demographics data for the general population of Germany; 891 211 persons (rows) x 366 features (columns).
- 2) Udacity\_CUSTOMERS\_052018.csv: Demographics data for customers of a mail-order company; 191 652 persons (rows) x 369 features (columns).
- 3) Udacity\_MAILOUT\_052018\_TRAIN.csv: Demographics data for individuals who were targets of a marketing campaign; 42 982 persons (rows) x 367 (columns).

4) Udacity\_MAILOUT\_052018\_TEST.csv: Demographics data for individuals who were targets of a marketing campaign; 42 833 persons (rows) x 366 (columns).

In addition, two description files have been provided to give attribute information:

5) DIAS Information Levels - Attributes 2017.xlsx: A top-level list of attributes and descriptions, organized by informational category.

6) DIAS Attributes - Values 2017.xlsx file: A detailed mapping of data values for each feature in alphabetical order.

### **Problem Statement**

The problem that we will be investigating can be formulated as:

*“From demographic data of a person, how can a mail order company acquire new customers in an efficient way?”*

This problem requires us to consider what we can do to predict with high accuracy whether a person with associated demographics data will be a new customer to our client company. Furthermore, how can we predict with confidence the probability of people with demographic profiles turning into future customers?

### **Evaluation Metrics**

The selected supervised machine learning model will be used to make predictions on the mailout campaign data. The evaluation metric is the Area Under the Receiver Operating Characteristics Curve (AUC for the ROC curve).

The decision of selecting the AUC-ROC curve to measure performance is due to the fact that it is one of the best performance evaluation techniques in imbalanced cases. In our case, the MAILOUT\_TRAIN dataset displays large class imbalance, where most individuals did not respond to the mailout (Fig.1). Thus, predicting individual classes and using accuracy does not seem to be an appropriate performance evaluation method.

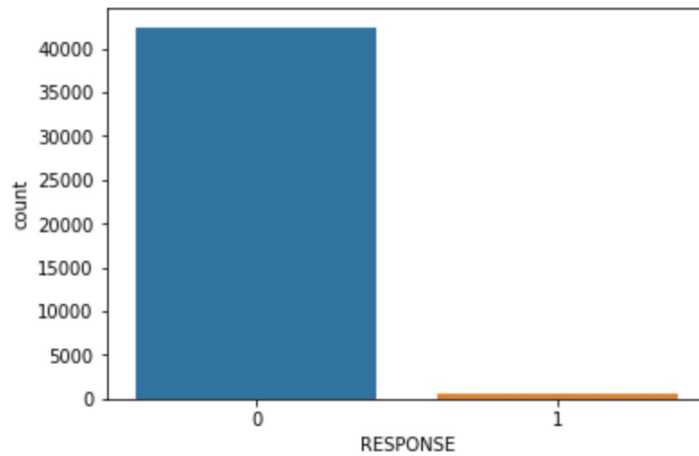


Figure 1. Counts of responses of individuals in the MAILOUT\_TRAIN dataset.

The interpretation of the AUC-ROC curve is available in the Proposal file in the GitHub repository.

## II. Analysis

### Data Exploration

Exploring the AZDIAS and CUSTOMERS datasets is the first step in this project. The two datasets contain the expected number of rows and columns as described. Each row of the two datasets represents a single person, but also includes information outside of individuals, including information about their household, building, and neighborhood. Compared with AZDIAS, CUSTOMERS contains three extra columns ('CUSTOMER\_GROUP', 'ONLINE\_PURCHASE', and 'PRODUCT\_GROUP'), which provide broad information about the customers depicted in the file.

The MAILOUT train and test datasets share similar structures with the two datasets mentioned above. However, the MAILOUT dataset includes one additional column, which indicated whether or not each recipient became a customer of the company. For the "TRAIN" subset, this column has been retained, but in the "TEST" subset it has been removed.

The two metadata file helps in understanding what each feature represents and cleaning unknown values.

## **Data Pre-processing**

Exploring data helps gain an insight into data points and features. However, since it is required to use complete data points to train a model, any missing or mis-recorded values requires to be cleaned. Here, we start with examining AZDIAS. The other three datasets will be pre-processed in a similar fashion.

### **1. Addressing unknown values**

The first step is to deal with unknown representations across all the columns. In the attributes-values spreadsheet, unknown records are represented by one or multiple numerical values. These unknown values should be treated as missing entries.

Thus, they were changed to NAN values.

### **2. Exploring missing values both column-wise and row-wise**

After replacing unknown values with NAN values, missing values were studied both column-wise and row-wise.

We investigated the proportion of NAN values per column. It turned out that there were many missing values in the AZDIAS, especially eight features having more than 60% missing values (Fig.2). To figure out a threshold for dropping columns, we plotted number of columns at different percentage of missing values. Since most of the columns had percentage of missing values of less than 20% (Fig.3), we decided to take 20% as the cut-off.

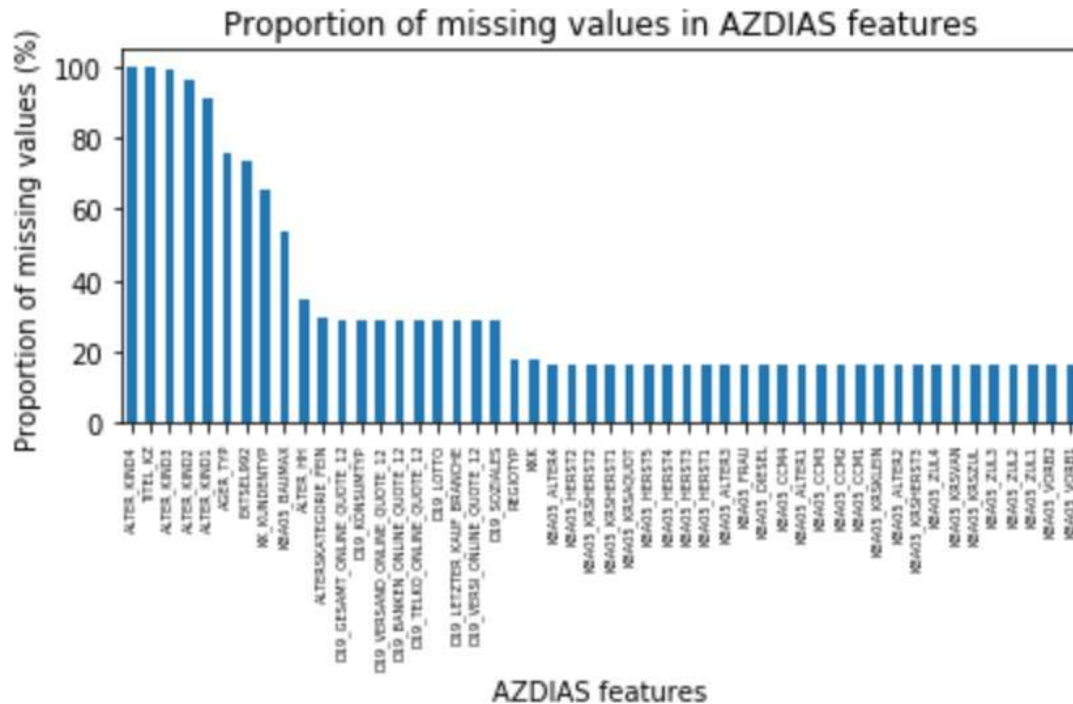


Figure 2. Percentage of missing values per ASDIAS features

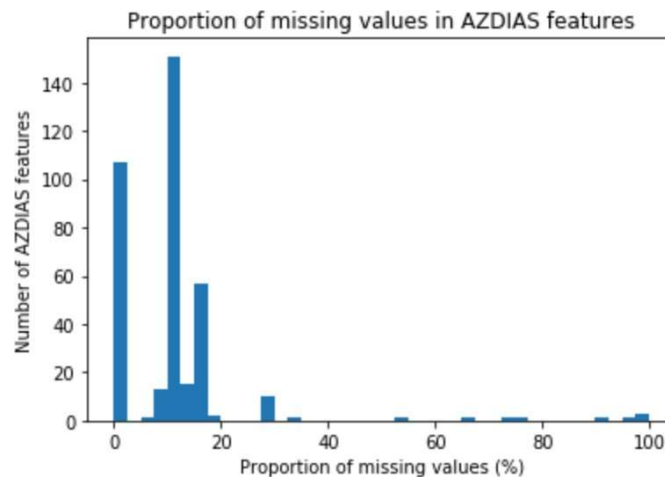


Figure 3. Number of columns for the proportion of missing values in AZDIAS.

We next explored missing values per row. Similarly, we were trying to figure out a threshold to drop rows containing too many missing values. Most of the columns had percentage of missing values of less than 10%. However, there are also a high number of rows with missing values greater than 60% (Fig.4). We decided to drop any rows that have over 10% missing values.

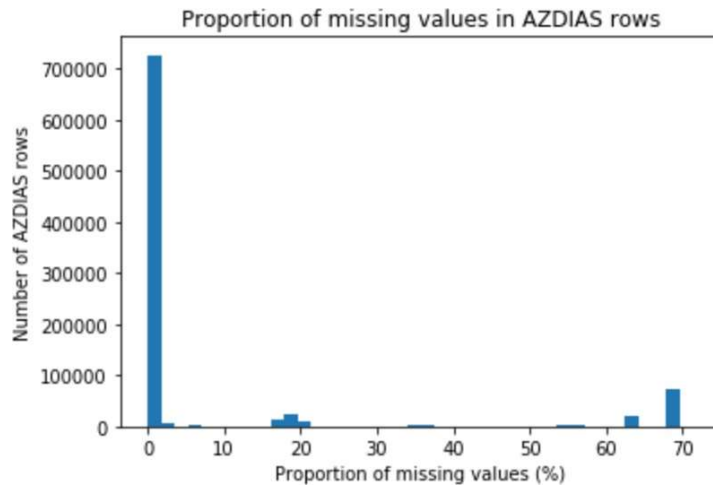


Figure 4. Number of rows for the proportion of missing values in AZDIAS.

### 3. Re-encoding categorical and mixed features

Some columns in the AZDIAS dataset contain categorical values. Some other may include mixed types. To prepare these for feature extraction, we'll want to convert these into numerical values.

During the exploration of categorical and mixed features, we managed to convert 'EINGEFUEGT\_AM' to a datetime type and treat 'X' and 'XX' as NAN values.

For binary categorical variables that have non-numerical values, we need to convert them to numerical data. Specifically, 'OST\_WEST\_KZ' was the binary categorical variable that was converted to numerical values based upon the DIAS Attributes - Values 2017.xlsx file.

For multi-level categorical variables, they may increase the number of features for modeling. Therefore, we may consider dropping them. We dropped 'CAMEO\_DE U\_2015' and 'EINGEFUEGT\_AM' in this step since both of them were with more than 10 values.

In terms of mixed features, we need to convert them to numerical values. There exist 4 mixed features, including 'CAMEO\_INTL\_2015', 'LP\_LEBENS PHASE\_FEIN', 'LP\_LEBENS PHASE\_GROB' and 'PRAEGENDE\_JUGENDJAHRE'. For these, we need to split them into multiple single-type features. 'CAMEO\_INTL\_2015' was re-encoded into 'CAMEO\_INTL\_2015\_WEALTH\_LEVEL' and 'CAMEO\_INTL\_2015\_STATUS'. Features that started with LP had duplicated information. A function was developed to handle these features. The 4th mixed-feature column was



'PRAEGENDE\_JUGENDJAHRE'. We re-encoded this feature to contain two values being Mainstream and Avantgarde, respectively.

#### **4. Dropping highly correlated features**

We will be implementing feature correlation to determine too high-correlated features since they many over-inflate the importance of a single feature. In this project, we defined too highly correlated features as having correlations with a column over 0.9. In this step, 15 features were dropped.

#### **5. Completing missing values**

After removing features and rows having proportion of missing values over the thresholds, there still exist missing values in the data. For the remaining missing values, we imputed them by using median values in each column.

#### **6. Scaling features**

Prior to applying PCA, we need to scale features to be of the same range since PCA may be influenced by variations in scales of features. A standard scaler was used to bring all the features to the same range.

### **III. Algorithms, Techniques and Methodology**

#### **Customer Segmentation**

The goal of the first part of the project is to characterize the relationship between existing customers and the demographics of the German population, and to find out which types of people in the general population are more likely to become customers.

#### **1. Dimensionality Reduction**

Before feeding data into a machine learning model, a dimensionality reduction step is a necessity with the aim being forming a smaller set of features to better help separate our data. The technique that was used to reduce the number of features was the principal component analysis (PCA).

The idea of reducing dimensionality is that we want to keep a small number of features while retaining high explained data variance. To decide how many top components to include, it's helpful to look at how much data variance the components capture. From the two figures below (Fig.5 and Fig.6), we can conclude that 160 features can explain around 85% data variance. Further check on explained data variance by 160 principal components demonstrated that they could explain 86% data variance in AZDIAS and 85% data variance in CUSTOMERS.

Therefore, I decided to keep 160 principal components since they could retain the main information in both datasets.

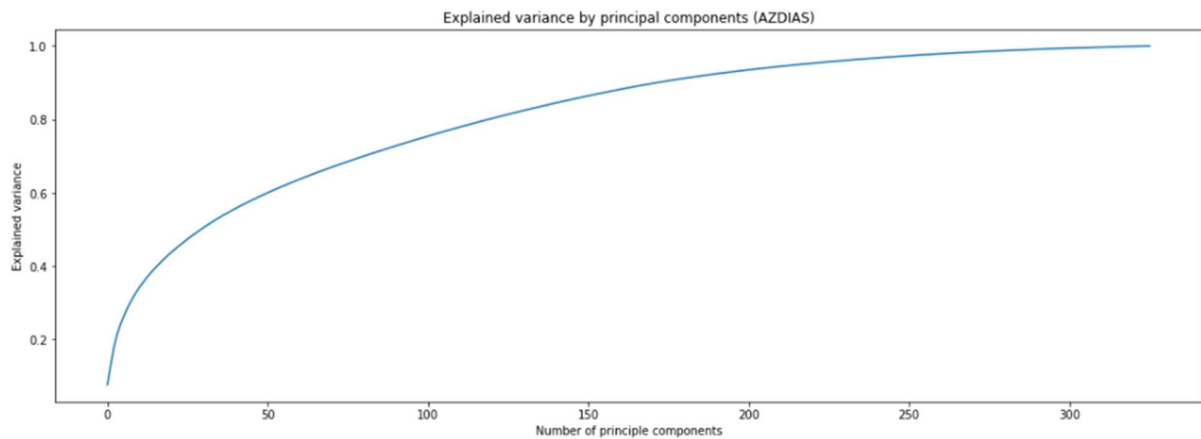


Figure 5. cumulative explained data variance by principal components in AZDIAS.

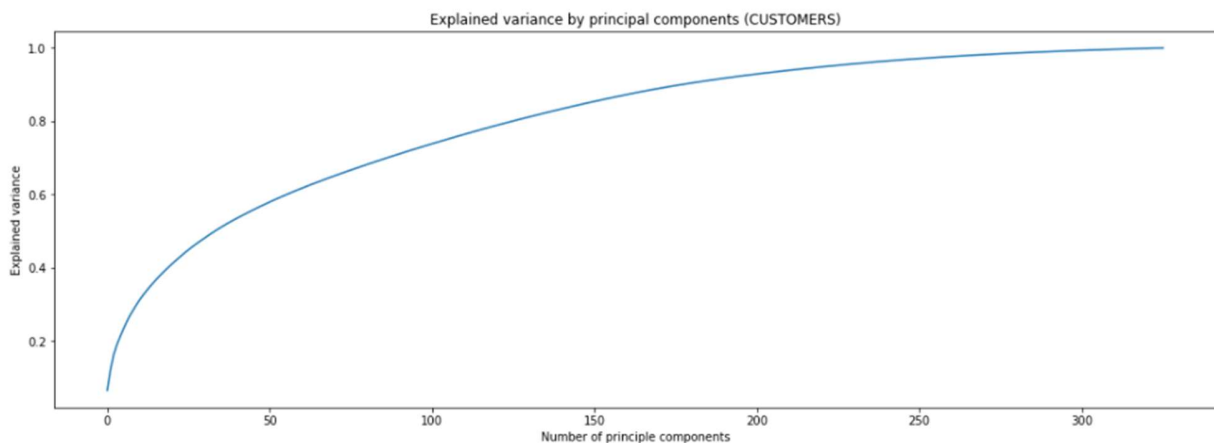


Figure 6. Cumulative explained data variance by principal components in CUSTOMERS.

## 2. Clustering

After having PCA attributes set up, the next is to divide the general population and the customer base into different segments. The unsupervised clustering algorithm, k-means, is chosen to segment the general population and customers. We want to

select a  $k$  such that data points in a single cluster are close together but that there are enough clusters to effectively separate the data. After trying several values for  $k$ , the centroid distance typically reaches some "elbow"; it stops decreasing at a sharp rate and this indicates a good value of  $k$ .

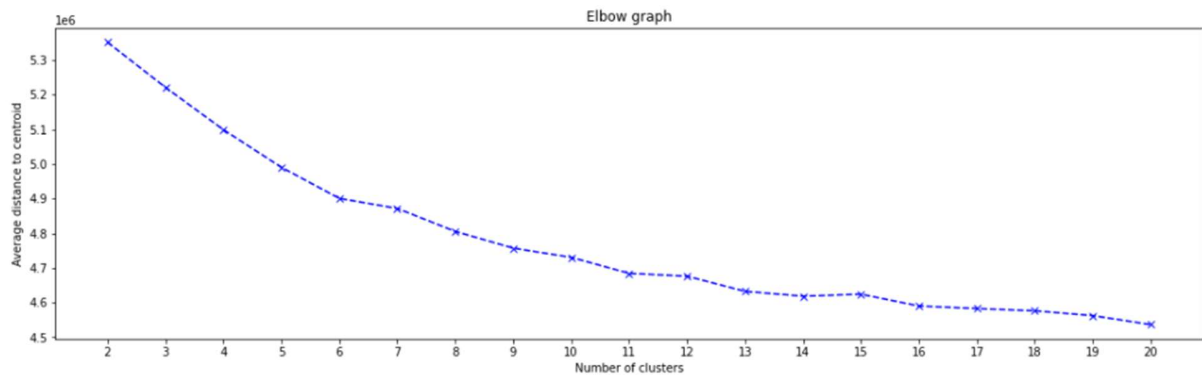


Figure 7. Elbow graph showing within-cluster distance for different number of clusters. From the elbow graph (Fig. 7), we can see that within-cluster distances decrease at a high rate until reaching 11 clusters. This indicates that there is enough separation to distinguish the data points in each cluster, but also that we included enough clusters so that the data points aren't extremely far away from each cluster. Thus, we came to the conclusion that 11 was the optimal number of clusters.

We next looked into clusters and compared the general population and customer base for each of the 11 clusters. We can know that some clusters were overrepresented by customers whereas others didn't (Fig. 8). This gives us an idea of better targeting future customers.

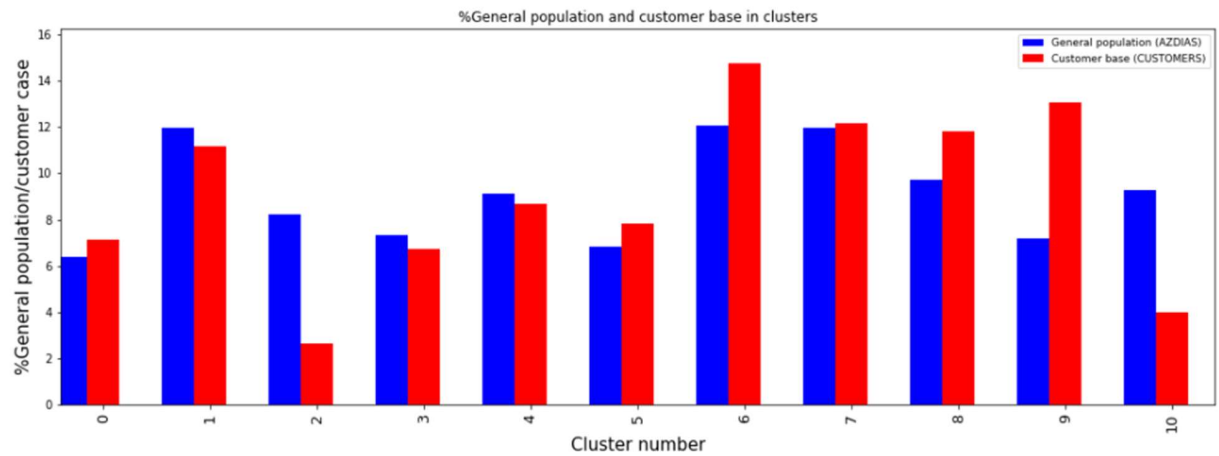


Figure 8. Percentage of general population vs customer base in each cluster.

### Customer Acquisition

The second part of this project is to use supervised learning techniques to predict potential customers based on demographic data. We will be using the MAILOUT train and test datasets. As mentioned, the MAILOUT dataset has one additional column included in the training dataset, which indicated whether or not a person will become a customer of the company. The two datasets were pre-processed and scaled using similar strategies applied to AZDIAS and CUSTOMERS. Additionally, in the MAILOUT\_TRAIN dataset, we observed large class imbalance (Fig.1). To deal with this, we resampled the minority class, leading to the same number of observations in both response classes.

We have tried the following supervised learning algorithms (each with performances):

#### 1) Benchmark Model: Logistic Regression

The first step in supervised learning is to set a benchmark, which has the base performance with a simple model. The benchmark model will be used to compare and evaluate performances of other trained models. Logistic Regressor is widely used as a benchmark model [4].

```
: logistic_reg = LogisticRegression(random_state=10)
```

```
: print(train_classifier(logistic_reg, {}))
```

```
Time taken: 15.302631616592407  
AUROC score: 0.7869742508616179  
LogisticRegression(random_state=10)
```

## 2) Random Forest Classifier

Random Forest Classifier takes fast runtime, and it is able to deal with unbalanced and missing data, which is suitable for our training dataset. However, it may overfit data that are particularly noisy [5].

```
random_forest_clf = RandomForestClassifier(random_state=10)
```

```
print(train_classifier(random_forest_clf, {}))
```

```
Time taken: 143.02382254600525  
AUROC score: 0.9935088703585709  
RandomForestClassifier(random_state=10)
```

## 3) Gradient Boosting Classifier

### **GradientBoostingClassifier**

```
gradient_boosting_clf = GradientBoostingClassifier(random_state=10)
```

```
print(train_classifier(gradient_boosting_clf, {}))
```

```
Time taken: 425.0988471508026  
AUROC score: 0.9024049565159767  
GradientBoostingClassifier(random_state=10)
```

Gradient boosting often provides high predictive accuracy and offers lots of flexibility [6].

#### 4) Light GBM Classifier

Light GBM is a gradient boosting framework that uses tree-based learning algorithm [7]. It is gaining extreme popularity because of it is fast, takes lower memory to run and has better accuracy than any other boosting algorithm. It could be a good complimentary supervised learning approach to be applied to our training dataset. However, it may grow a deeper decision tree, being prone to over-fitting [7].

##### LightGBMClassifier

Light GBM can handle the large size of data and takes lower memory to run. Another reason why Light GBM is so popular is because it focuses on accuracy of results.

```
j: lightgbm_clf = LGBMClassifier(random_state=10)
```

```
j: print(train_classifier(lightgbm_clf, {}))
```

```
Time taken: 15.292855739593506  
AUROC score: 0.9874326898738097  
LGBMClassifier(random_state=10)
```

#### 5) Ada Boost Classifier

Adaptive Boosting is one of ensemble boosting classifier. It combines multiple classifiers to increase the accuracy of classifiers. AdaBoost is easy to implement and not prone to over-fitting. The disadvantage of this approach is that it is highly affected by outliers because it tries to fit each data point perfectly [8]. Since we didn't look deeply on outliers in the data, AdaBoost may not perform very well.

##### AdaBoostClassifier

```
: adaboost_clf = AdaBoostClassifier(random_state=10)
```

```
: print(train_classifier(adaboost_clf, {}))
```

```
Time taken: 98.37448906898499  
AUROC score: 0.818519529664114  
AdaBoostClassifier(random_state=10)
```

#### 6) XG Boost Classifier

```
Time taken: 65.22626829147339  
AUROC score: 0.99258412283386  
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,  
               colsample_bynode=1, colsample_bytree=1, enable_categorical=False,  
               gamma=0, gpu_id=-1, importance_type=None,  
               interaction_constraints='', learning_rate=0.300000012,  
               max_delta_step=0, max_depth=6, min_child_weight=1, missing=nan,  
               monotone_constraints=(), n_estimators=100, n_jobs=12,  
               num_parallel_tree=1, predictor='auto', random_state=10,  
               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,  
               tree_method='exact', validate_parameters=1, verbosity=None)
```

XGBoost stands for Extreme Gradient Boosting, which is an improved version of the Gradient Boosting Algorithm [9].

## IV. Results

After hyperparameter tuning, the best trained model was the Gradient Boosting Classifier algorithm with an AUROC score being 0.902, which displayed better performance compared with the benchmark model.

### Fine tuning GradientBoostingClassifier

```
: gbc = GradientBoostingClassifier(random_state=10)

gbc_param_grid = {'loss': ['deviance'],
                  'max_depth': [3],
                  'learning_rate': [0.1],
                  'n_estimators': [100],
                  'random_state': [10]
                  }

gbc_best_estimator = train_classifier(gbc, gbc_param_grid)

Time taken: 419.3621542453766
AUROC score: 0.9024049565159767
```

Final predictions were made on the preprocessed and scaled MAILOUT\_TEST dataset. Among 42,833 individuals listed in the test dataset, 9,115 individuals are likely to become a customer of the mail-order company, which accounts for approximately 21% of the test individuals.

## **V. Justification**

The logistic Regressor achieved an AUROC score of 0.786 on the training dataset MAILOUT\_TRAIN, whereas the Gradient Boosting Classifier had scored 0.90.

This demonstrated that the selected Gradient Boosting Classifier model had surpassed the benchmark model on the performances on the training dataset, which make our chosen model potentially a good candidate for the final supervised model.

Light GBM classifier, XG Boost Classifier and Random Forest Classifier had AUROC scores being close to 1, which raised concerns on overfitting issues.

## **VI. Improvements and Future Steps**

Although I got a very high scores, there is still room for improvement.

We can consider including the following steps in future improvements:

1. Addressing more multi-level categorical features and one-hot encoding them.
2. Look deeper to find and understand outliers.
3. Dealing with class imbalance using down-sampling methods instead of oversampling.
4. Setting different ranges of hyperparameters in the fine tuning steps.
5. Pick less PCA components than 160: ranging from 130 to 160 for example.



## References

- [1] Arvato. *Wikipedia*. <https://en.wikipedia.org/wiki/Arvato>
- [2] Customer Segmentation. *Wikipedia*.  
[https://en.wikipedia.org/wiki/Market\\_segmentation](https://en.wikipedia.org/wiki/Market_segmentation)
- [3] Udacity+Arvato: Identify Customer Segments. *Kaggle*.  
<https://www.kaggle.com/c/udacity-arvato-identify-customers/overview/evaluation>
- [4] LOGISTIC REGRESSION CLASSIFIER. *Towards Data Science*.  
<https://towardsdatascience.com/logistic-regression-classifier-8583e0c3cf9>
- [5] Random Forest. *Wikipedia*. [https://en.wikipedia.org/wiki/Random\\_forest](https://en.wikipedia.org/wiki/Random_forest)
- [6] Understanding Gradient Boosting Machines. *Towards Data Science*.  
<https://towardsdatascience.com/understanding-gradient-boosting-machines-9be756fe76ab>
- [7] LightGBM for machine learning. *Programmer sought*.  
<https://www.programmersought.com/article/52844993233/>
- [8] AdaBoost Classifier in python. *DataCamp*.  
<https://www.datacamp.com/community/tutorials/adaboost-classifier-python>
- [9] LightGBM vs XGBoost – Which algorithm is better.  
<https://www.geeksforgeeks.org/lightgbm-vs-xgboost-which-algorithm-is-better/>
- [10] Sensitivity and Specificity. In *Wikipedia*. Retrieved from:  
[https://en.wikipedia.org/wiki/Sensitivity\\_and\\_specificity](https://en.wikipedia.org/wiki/Sensitivity_and_specificity)
- [11] Classification: ROC Curve and AUC. In *Google Developers*. Retrieved from: <https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc>