

# Arvato Financial Services, Customer Segmentation Project Report

## Table of Contents

<b>Project Description and Introduction:</b> .....	<b>1</b>
<b>Data Wrangling and Processing:</b> .....	<b>2</b>
<b>Customer Segmentation:</b> .....	<b>4</b>
Reducing Dimensions: .....	4
Creating Clusters: .....	5
Cluster Comparison: .....	5
<b>Supervised Learning and Output Predictions</b> .....	<b>6</b>
Analyze data and split the input features and label .....	6
Model Evaluation and Prediction: .....	7
Popular Features: .....	7
<b>Conclusion</b> .....	<b>8</b>
<b>Things I would like to do:</b> .....	<b>9</b>
<b>References:</b> .....	<b>9</b>

## Project Description and Introduction:

**Domain Background:** Financial Services, Arvato is an IT company that has different products lines. We will be focusing on a use case that will help their financial group to identify potential customers through their mail order campaign.

**Problem Statement:** In this project we will analyze demographics data of customers of a mail-order sales company (Arvato) in Germany, comparing it against demographics information for the general population. We will use unsupervised learning techniques to perform customer segmentation and identifying parts of the population that best describe the core customer base of the company. Then, we'll apply what we've learned on a dataset, with demographics information for targets of a marketing campaign for the company, and use a model to predict which individuals are most likely to convert into becoming customers for the company.

Due to the size and imbalance of the dataset the evaluation metric of AUC for the ROC curve, where ROC curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) will be used to evaluate the model.

### Key Project Milestones:

- Customer Segmentation using Demographics Data
- Develop a Model to predict potential customer for the company through ad campaign
- Upload the results to Kaggle competition to share and score the performance and analysis

# Data Wrangling and Processing:

As with any data science or machine learning project, data cleanup and processing is the most crucial and time consuming aspect of the project. The data is what makes or breaks a model.

For our project, we have been give four real-life datasets of population, customer data for Germany for the mail order company and the user data of the mail order campaign.

- **AZDIAS:** Demographics data for the general population of Germany
- **CUSTOMERS:** Demographics data for customers of a mail-order company
- **MAILOUT\_TRAIN** and **MAILOUT\_TEST:** Demographics data for individuals, targets of a marketing campaign

First step in the project is the data analysis. We will be analyzing the data and process it which will help us identify clusters for our customer base and will eventually help us identify potential customers for the company.

Initial sneak peek at the demographics data AZDIAS which contains 891 211 persons (rows) x 366 features (columns). It has a combination of numeric and categorical data. The first part of the process is to identify all the NaNs in the dataset and handle it to make the dataset conducive for the unsupervised learning part of the project.

Pandas Profiling Report

Overview

Variables

Correlations

Missing values

Sample

Dataset info

Number of variables

366

Number of observations

5

Missing cells

284 (15.5%)

Duplicate rows

0 (0.0%)

Total size in memory

14.4 KiB

Average record size in memory

2.9 KiB

Variables types

Numeric

1

Categorical

24

Boolean

14

Date

0

URL

0

Text (Unique)

0

Rejected

327

Unsupported

0

Warnings

AKT\_DAT\_KL has 1 (20.0%) missing values

Missing

ALTER\_HH has 1 (20.0%) missing values

Missing

ALTER\_KIND1 has constant value "nan"

Rejected

ALTER\_KIND2 has constant value "nan"

Rejected

ALTER\_KIND3 has constant value "nan"

Rejected

## Exploring and Analysing the Dataset

Step 1: Sneak Peak at the data

In [4]:	population_data.describe()																																																																																	
Out[4]:	<table><tr><th></th><th>LNR</th><th>AGER_TYP</th><th>AKT_DAT_KL</th><th>ALTER_HH</th><th>ALTER_KIND1</th><th>ALTER_KIND2</th><th>ALTER_KIND3</th><th>ALTER_K</th></tr><tr><td>count</td><td>8.912210e+05</td><td>891221.000000</td><td>817722.000000</td><td>817722.000000</td><td>81058.000000</td><td>29499.000000</td><td>6170.000000</td><td>1205.000</td></tr><tr><td>mean</td><td>6.372630e+05</td><td>-0.358435</td><td>4.421928</td><td>10.864126</td><td>11.745392</td><td>13.402658</td><td>14.476013</td><td>15.08962</td></tr><tr><td>std</td><td>2.572735e+05</td><td>1.198724</td><td>3.638805</td><td>7.639683</td><td>4.097660</td><td>3.243300</td><td>2.712427</td><td>2.452932</td></tr><tr><td>min</td><td>1.916530e+05</td><td>-1.000000</td><td>1.000000</td><td>0.000000</td><td>2.000000</td><td>2.000000</td><td>4.000000</td><td>7.000000</td></tr><tr><td>25%</td><td>4.144580e+05</td><td>-1.000000</td><td>1.000000</td><td>0.000000</td><td>8.000000</td><td>11.000000</td><td>13.000000</td><td>14.00000</td></tr><tr><td>50%</td><td>6.372630e+05</td><td>-1.000000</td><td>3.000000</td><td>13.000000</td><td>12.000000</td><td>14.000000</td><td>15.000000</td><td>15.00000</td></tr><tr><td>75%</td><td>8.600680e+05</td><td>-1.000000</td><td>9.000000</td><td>17.000000</td><td>15.000000</td><td>16.000000</td><td>17.000000</td><td>17.00000</td></tr><tr><td>max</td><td>1.082873e+06</td><td>3.000000</td><td>9.000000</td><td>21.000000</td><td>18.000000</td><td>18.000000</td><td>18.000000</td><td>18.00000</td></tr></table>		LNR	AGER_TYP	AKT_DAT_KL	ALTER_HH	ALTER_KIND1	ALTER_KIND2	ALTER_KIND3	ALTER_K	count	8.912210e+05	891221.000000	817722.000000	817722.000000	81058.000000	29499.000000	6170.000000	1205.000	mean	6.372630e+05	-0.358435	4.421928	10.864126	11.745392	13.402658	14.476013	15.08962	std	2.572735e+05	1.198724	3.638805	7.639683	4.097660	3.243300	2.712427	2.452932	min	1.916530e+05	-1.000000	1.000000	0.000000	2.000000	2.000000	4.000000	7.000000	25%	4.144580e+05	-1.000000	1.000000	0.000000	8.000000	11.000000	13.000000	14.00000	50%	6.372630e+05	-1.000000	3.000000	13.000000	12.000000	14.000000	15.000000	15.00000	75%	8.600680e+05	-1.000000	9.000000	17.000000	15.000000	16.000000	17.000000	17.00000	max	1.082873e+06	3.000000	9.000000	21.000000	18.000000	18.000000	18.000000	18.00000
	LNR	AGER_TYP	AKT_DAT_KL	ALTER_HH	ALTER_KIND1	ALTER_KIND2	ALTER_KIND3	ALTER_K																																																																										
count	8.912210e+05	891221.000000	817722.000000	817722.000000	81058.000000	29499.000000	6170.000000	1205.000																																																																										
mean	6.372630e+05	-0.358435	4.421928	10.864126	11.745392	13.402658	14.476013	15.08962																																																																										
std	2.572735e+05	1.198724	3.638805	7.639683	4.097660	3.243300	2.712427	2.452932																																																																										
min	1.916530e+05	-1.000000	1.000000	0.000000	2.000000	2.000000	4.000000	7.000000																																																																										
25%	4.144580e+05	-1.000000	1.000000	0.000000	8.000000	11.000000	13.000000	14.00000																																																																										
50%	6.372630e+05	-1.000000	3.000000	13.000000	12.000000	14.000000	15.000000	15.00000																																																																										
75%	8.600680e+05	-1.000000	9.000000	17.000000	15.000000	16.000000	17.000000	17.00000																																																																										
max	1.082873e+06	3.000000	9.000000	21.000000	18.000000	18.000000	18.000000	18.00000																																																																										
	8 rows x 360 columns																																																																																	

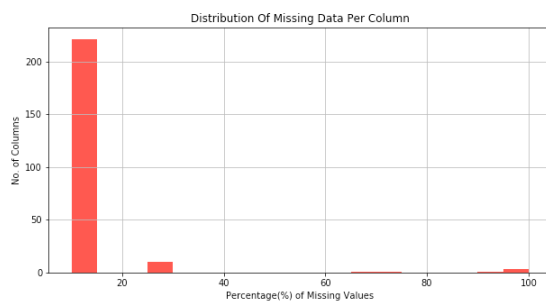
In [5]:	customers_data.describe()																																																																																	
Out[5]:	<table><tr><th></th><th>LNR</th><th>AGER_TYP</th><th>AKT_DAT_KL</th><th>ALTER_HH</th><th>ALTER_KIND1</th><th>ALTER_KIND2</th><th>ALTER_KIND3</th><th>ALTER_K</th></tr><tr><td>count</td><td>191652.000000</td><td>191652.000000</td><td>145056.000000</td><td>145056.000000</td><td>11766.000000</td><td>5100.000000</td><td>1275.000000</td><td>236.000</td></tr><tr><td>mean</td><td>95826.500000</td><td>0.344359</td><td>1.747525</td><td>11.352009</td><td>12.337243</td><td>13.672353</td><td>14.647059</td><td>15.3771</td></tr><tr><td>std</td><td>55325.311233</td><td>1.391672</td><td>1.966334</td><td>6.275026</td><td>4.006050</td><td>3.243335</td><td>2.753787</td><td>2.30765</td></tr><tr><td>min</td><td>1.000000</td><td>-1.000000</td><td>1.000000</td><td>0.000000</td><td>2.000000</td><td>2.000000</td><td>5.000000</td><td>8.00000</td></tr><tr><td>25%</td><td>47913.750000</td><td>-1.000000</td><td>1.000000</td><td>8.000000</td><td>9.000000</td><td>11.000000</td><td>13.000000</td><td>14.0000</td></tr><tr><td>50%</td><td>95826.500000</td><td>0.000000</td><td>1.000000</td><td>11.000000</td><td>13.000000</td><td>14.000000</td><td>15.000000</td><td>16.0000</td></tr><tr><td>75%</td><td>143739.250000</td><td>2.000000</td><td>1.000000</td><td>16.000000</td><td>16.000000</td><td>16.000000</td><td>17.000000</td><td>17.0000</td></tr><tr><td>max</td><td>191652.000000</td><td>3.000000</td><td>9.000000</td><td>21.000000</td><td>18.000000</td><td>18.000000</td><td>18.000000</td><td>18.0000</td></tr></table>		LNR	AGER_TYP	AKT_DAT_KL	ALTER_HH	ALTER_KIND1	ALTER_KIND2	ALTER_KIND3	ALTER_K	count	191652.000000	191652.000000	145056.000000	145056.000000	11766.000000	5100.000000	1275.000000	236.000	mean	95826.500000	0.344359	1.747525	11.352009	12.337243	13.672353	14.647059	15.3771	std	55325.311233	1.391672	1.966334	6.275026	4.006050	3.243335	2.753787	2.30765	min	1.000000	-1.000000	1.000000	0.000000	2.000000	2.000000	5.000000	8.00000	25%	47913.750000	-1.000000	1.000000	8.000000	9.000000	11.000000	13.000000	14.0000	50%	95826.500000	0.000000	1.000000	11.000000	13.000000	14.000000	15.000000	16.0000	75%	143739.250000	2.000000	1.000000	16.000000	16.000000	16.000000	17.000000	17.0000	max	191652.000000	3.000000	9.000000	21.000000	18.000000	18.000000	18.000000	18.0000
	LNR	AGER_TYP	AKT_DAT_KL	ALTER_HH	ALTER_KIND1	ALTER_KIND2	ALTER_KIND3	ALTER_K																																																																										
count	191652.000000	191652.000000	145056.000000	145056.000000	11766.000000	5100.000000	1275.000000	236.000																																																																										
mean	95826.500000	0.344359	1.747525	11.352009	12.337243	13.672353	14.647059	15.3771																																																																										
std	55325.311233	1.391672	1.966334	6.275026	4.006050	3.243335	2.753787	2.30765																																																																										
min	1.000000	-1.000000	1.000000	0.000000	2.000000	2.000000	5.000000	8.00000																																																																										
25%	47913.750000	-1.000000	1.000000	8.000000	9.000000	11.000000	13.000000	14.0000																																																																										
50%	95826.500000	0.000000	1.000000	11.000000	13.000000	14.000000	15.000000	16.0000																																																																										
75%	143739.250000	2.000000	1.000000	16.000000	16.000000	16.000000	17.000000	17.0000																																																																										
max	191652.000000	3.000000	9.000000	21.000000	18.000000	18.000000	18.000000	18.0000																																																																										
	8 rows x 361 columns																																																																																	

I followed a series of data cleaning steps to prepare the data for the next stages of the project. Some of them were like the following

- Identified NaNs or the unknowns from the dataset over rows and columns
- Removed columns with high correlation to reduce the feature set > 65%
- Removed rows with most unknowns or NaNs by identifying threshold < 16%
- Encoded the data for the remaining of the NaNs or unknowns in the dataset
- Handle the outliers in the dataset
- Scaled and standardized the dataset, used Mode as the filler value

```
In [7]: #Top 10 Candidates
population_data_NaN_Percent.sort_values(ascending=False).head(10)

Out[7]: ALTER_KIND4          99.864792
ALTER_KIND3          99.307691
ALTER_KIND2          96.690047
ALTER_KIND1          90.904837
EXTSEL992           73.399639
KK_KUNDENTYP         65.596749
ALTERSKATEGORIE_FEIN 29.504130
D19_LETZTER_KAUF_BRANCHE 28.849522
D19_LOTTO            28.849522
D19_VERSI_ONLINE_QUOTE_12 28.849522
dtype: float64
```



#### Step 3: Remove NaNs from the Dataset

```
In [12]: #Removing all rows from the dataset with threshold less than 16
population_data = population_data[population_data.isnull().sum(axis=1) <= 16].reset_index(drop=True)
print('New Dataset Size After Removing the Rows: ', population_data.shape[0])

New Dataset Size After Removing the Rows: 733227

In [13]: #Removing all columns from the dataset with threshold > 65%
Cols_To_Be_Dropped = population_data.columns[Column_Most_NaN > 0.65]
print('Columns with NaNs > 65%: ', Cols_To_Be_Dropped)

Columns with NaNs > 65%: Index(['ALTER_KIND1', 'ALTER_KIND2', 'ALTER_KIND3', 'ALTER_KIND4', 'EXTSEL992',
                                'KK_KUNDENTYP'],
                                dtype='object')
```

The dataset after all data preprocessing and encoding ended up looking like this, which will be used further.

```
In [44]: #do a check of missing values
population_data.head()
```

```
Out[44]:
```

	0	1	2	3	4	5	6	7	8	9	...	394	395	396	397	398	399	400	401	402	403
1	910225.0	-1.0	9.0	17.0	17.0	10.0	0.0	0.0	1.0	7.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
4	910244.0	3.0	1.0	10.0	10.0	5.0	0.0	0.0	1.0	2.0	...	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
5	910248.0	-1.0	9.0	0.0	9.0	4.0	0.0	0.0	1.0	3.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
6	910261.0	-1.0	1.0	14.0	14.0	6.0	0.0	0.0	1.0	5.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
9	645165.0	0.0	1.0	10.0	10.0	6.0	0.0	0.0	1.0	6.0	...	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	1.0

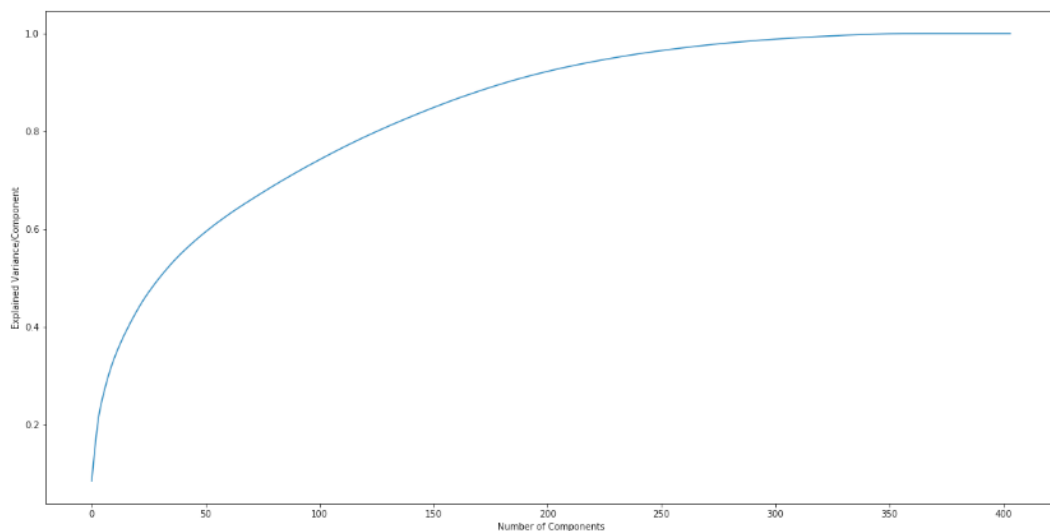
5 rows x 404 columns

## Customer Segmentation:

In this section of the project we will be using unsupervised learning and dimension reduction techniques to reduce the number of features and will enable to find the customer cluster and attributes from the general population.

### Reducing Dimensions:

Principal Component Analysis (PCA) is one of the most useful techniques in Exploratory Data Analysis to understand the data, reduce dimensions of data and for unsupervised learning in general. We will decide the number of input features to retain based on the cumulative variance.

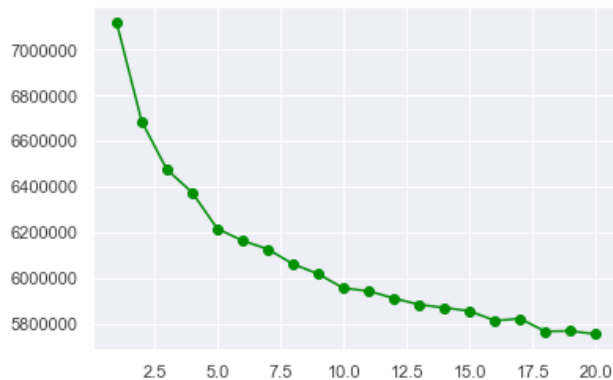


From the graph 260 is where the curve flattens and that's the number of principal components I have selected for the data with the cumulative variance around 95%. This the initial analysis, we can always come back and tune the numbers.

## Creating Clusters:

The next step in the process is to use the reduced features to identify clusters from the dataset. I have used the K-means algorithm to identify the number of clusters in the dataset and have used the same technique with the customer dataset to identify potential customer base.

I have used the elbow method to identify optimal number of clusters available in the dataset.

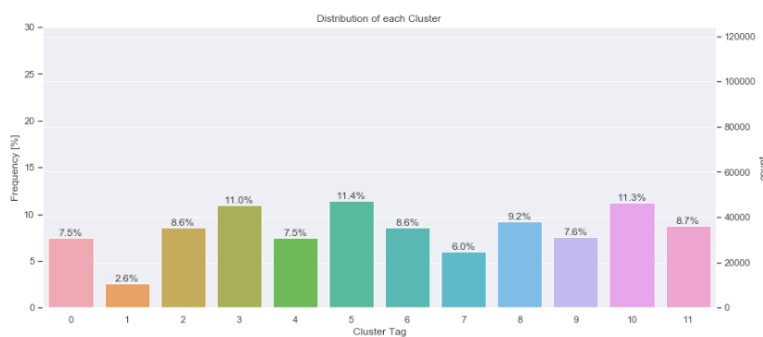


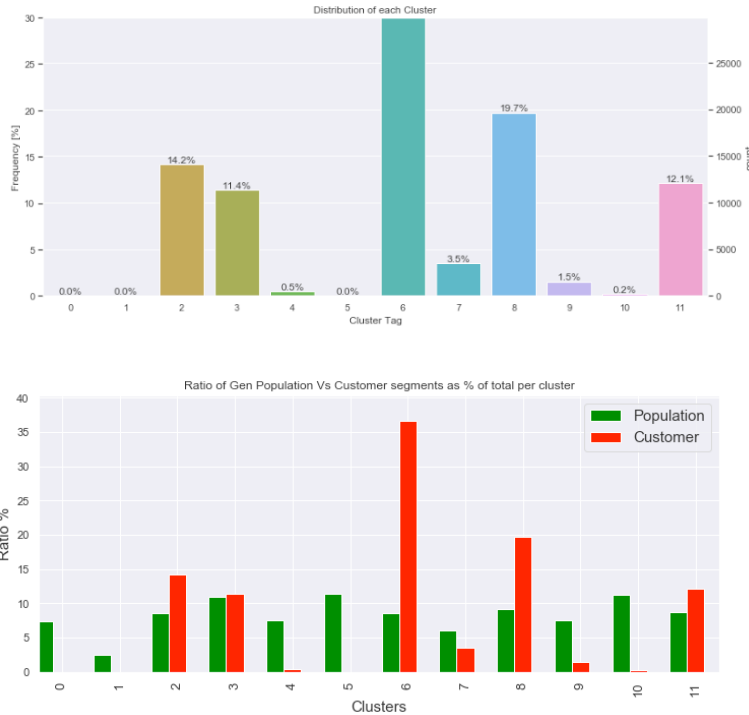
From the graph the elbow shape seems to form around 12 and after that the loss changes drastically. The ideal number of clusters will be 12.

I used K-means and the number of clusters identified above to fit the customer data to identify the potential base.

## Cluster Comparison:

In this section, I used the cluster data based on demographics of the general population of Germany and mapped the customer data from the mail-order sales company to those demographic clusters. The comparison of the two cluster distributions will enable us to identify the potential customer base.





Here we have two graphs that represents the distribution of population and customers of the company respectively. We can see the frequencies of cluster 6 and 8 in Customers data set are remarkably larger than the frequencies in Population (or general) dataset.

Considering the proportion of persons in each cluster for the general population and the customers the ratio should be fairly close but the mismatch of clusters between Population and Customers dataset indicates that there are only particular segments of the population that are interested in the company's products.

The huge difference of cluster 6/8 between two datasets suggests two things

- People in these two clusters could be potential target audience for the company
- Such large cluster difference could suggest the target base could have customers outside of the target demographics

## Supervised Learning and Output Predictions

In the last section of the project, we will apply supervised learning concepts to analyze the MAILOUT\_TRAIN and MAILOUT\_TEST dataset provided by the mail order company to predict probability of a person becoming a potential customer of the company.

### Analyze data and split the input features and label

I have split the dataset to identify the input features and label. The last column of the dataset is the label, RESPONSE we got from the potential customer.

On quick analysis, we could find that the dataset is not balanced as for 43k ads only 532 people responded to it. Due to the imbalance nature of the dataset I have made for the following decisions to handle the same

- Use the data cleanup function created earlier to clean the data
- Use stratified K-Fold technique to split the data for testing and training
- Use XGBoostRegressor algorithm as the problem is a classic classification problem
- Use roc\_auc\_score as a metric for evaluating the model

### Model Evaluation and Prediction:

There are really good algorithms like AdaBoostRegressor, GradientBoostingRegressor but I went with the XGBRegressor which works well for classification problems. I used 5-fold validation to train the data. I used techniques like GridSearchCV to identify the hyper parameters for the model.

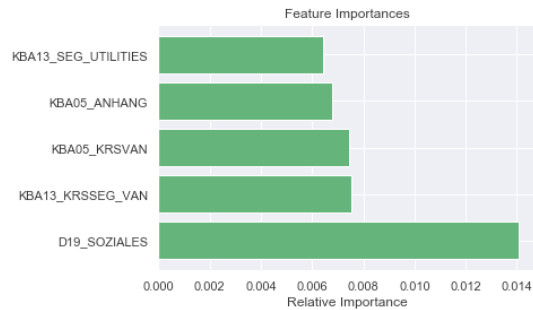
With the final set of parameters, I was able to achieve a roc\_auc\_score score of 0.78375.

#### Step 4: Finalize the parameters for algorithm

```
In [106]: clf_final = xgb.XGBRegressor(  
    objective = 'binary:logistic',  
    max_depth = 5,  
    min_child_weight = 2,  
    learning_rate = 0.01,  
    reg_alpha = 0.05,  
    subsample = 0.6,  
    colsample_bytree = 0.7,  
    gamma = 0.2,  
    scale_pos_weight = 1,  
    random_state = 42,  
    n_estimators = 500,  
    scoring = 'roc_auc')  
  
preds = clf_final.fit(X_train,y_train)  
predictions_test = preds.predict(X_val)  
print(roc_auc_score(y_val, predictions_test))  
  
0.7837570426762837
```

### Popular Features:

Popular Features identified from the dataset. D19\_SOZIALES turned out to be the most important features of all in the dataset. D19 series in the dataset corresponds to the transactional level details.



Using the same steps above, I have handled the Udacity\_MAILOUT\_052018\_TEST.csv dataset to predict the output for the Kaggle competition and exported the output in a csv file for submission.

**Step 3: Predict the Output**

```
In [160]: prediction_for_kaggle = clf_final.predict(mailout_test_clean)
```

```
In [161]: df_kaggle = pd.DataFrame(index=mailout_test['LNR'].astype('int32'), data=prediction_for_kaggle)
df_kaggle.rename(columns={0: "RESPONSE"}, inplace=True)
```

```
In [162]: df_kaggle.head(10)
```

Out[162]:

	RESPONSE
LNR	
1754	0.035161
1770	0.035284
1465	0.007356
1470	0.008121
1478	0.011761
1782	0.006371
1485	0.008821
1519	0.022835
1835	0.025162
1522	0.008129

**Step 4: Save the results**

```
In [163]: df_kaggle.to_csv('submission.csv')
```

## Conclusion

The overall project is a really good practical example of how to process real-life data and handle challenges while processing data and solve real life problems. The scope of the project ranges from data analysis to unsupervised learning to supervised learning along the way providing an ability to make decisions to achieve the end goal.

On a high level I have performed the following actions to solve the problem of identifying potential customers to the mail order company.

High Level Process Flow:

- Collect Data for the use case
- Clean up and Pre-process the demographics (General & Customer) data
- Feature Engineering
- PCA, reduce input features to identify effective input features
- Split the dataset into test, train and validation
- Generate the shoulder curve to identify the cluster size
- Apply K-means algorithm to the processed dataset to identify clusters



- Identify potential customer base clusters by comparing general vs customer cluster
- Clean up and Pre-process the given train and test customer data
- Identify Input Features
- Create a model
- Train the model
- Test the model with test customer data
- Validate output data and model performance
- Infer insights from the final predicted output

### Things I would like to do:

- Implement a similar solution in Sagemaker by using its built-in tools and available algorithms
- Try a different way to present the outcome and the results in a different way
- Would like to try different algorithms and train them for longer iteration to test the data
- Would like to implement deep learning model for the business problem

### References:

- <https://jakevdp.github.io/PythonDataScienceHandbook/index.html>
- <https://stackoverflow.com/questions/29294983/how-to-calculate-correlation-between-all-columns-and-remove-highly-correlated-on>
- <https://stats.stackexchange.com/questions/143700/which-is-better-replacement-by-mean-and-replacement-by-median>
- <https://towardsdatascience.com/ways-to-detect-and-remove-the-outliers-404d16608dba>
- <https://stackoverflow.com/questions/23199796/detect-and-exclude-outliers-in-pandas-data-frame>
- <https://machinelearningmastery.com/rescaling-data-for-machine-learning-in-python-with-scikit-learn/>