

Customer Segmentation Report for Arvato Financial Solutions

Project Overview

The objective of this project is to predict whether targets of a marketing campaign by a mail-order company will respond positively and become a customer. The project consists of two parts: unsupervised learning and supervised learning. In the unsupervised learning part, principal component analysis and KMeans clustering are used to determine the clusters that would likely describe the attributes of a likely customer. In the supervised learning part, a model utilising ensemble classifiers would be trained on the response of a marketing campaign to predict the response of a target customer.

Problem Statement

In this project, we will analyse the demographics data for customers of a mail-order sales company in Germany, comparing it against demographics information for the general population. The project is divided into 2 parts - unsupervised learning and supervised learning.

In the first part, unsupervised learning techniques shall be used to perform customer segmentation to identify the parts of the population that best describe the core customer base of the company.

In the second part, the attributes that best identifies the customer of the company as determined in the first part of the project will be applied to a third dataset with demographics information for targets of a marketing campaign for the company. A model will be used to predict which individuals are most likely to convert into becoming customers for the company.

Metrics

The metrics that will be used to evaluate our model is the 'area under the receiver operating characteristics curve' or known as the 'ROC_AUC' score. This metric is appropriate to evaluate our model as our dataset is highly imbalanced - there is a very low occurrence of true positives. The 'ROC_AUC' score represents the

probability that a classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one. It is a performance measurement for the classification problems at various threshold settings. ROC is a probability curve while AUC represents the degree of separability. It implies the ability of the model in distinguishing between classes. Basically, the higher the value of ROC_AUC, the better the model is at predicting 0s as 0s and 1s as 1s. Applied to our problem, this means the 'ROC_AUC' score represents the probability that a classifier will rank a random target that will respond as customer (1) than a random target that will not respond as non-customer(0).

Results, analysis & discussion

Data Provided

There are 4 datasets and 2 informational files.

The datasets are:

- 1.) Udacity_AZDIAS_052018.csv
- 2.) Udacity_CUSTOMERS_052018.csv
- 3.) Udacity_MAILOUT_052018_TEST.csv
- 4.) Udacity_MAILOUT_052018_TRAIN.csv

The information files are:

- 1.) DIAS Information Levels - Attributes 2017
- 2.) DIAS Attributes - Values 2017

The datasets generally represents attributes of each individual for each row.

Data Pre-processing

- 1.) Assess missing data in columns and rows

We start by assessing the dataset of Udacity_AZDIAS_052018.csv. To start assessing the missing data in columns, columns with missing values are identified. There are 2 types of missing data in the columns of the azdias dataset. The first type of missing values are values in columns with key codes that corresponds to 'unknown' or 'no_transaction_known'. The second type of missing values in columns occur simply because no values has been entered at all, it has been left blank.

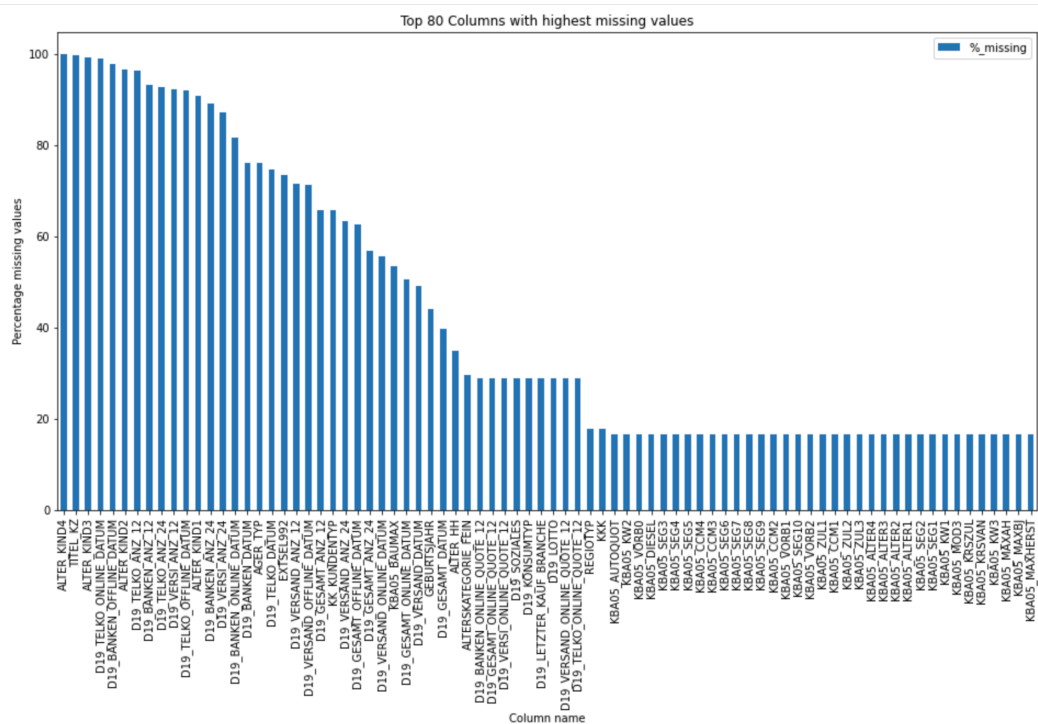
The first type of missing values in columns is addressed first. We iterate through the 'DIAS Attributes - Values 2017' file to identify attributes with key codes that corresponds to 'unknown' or 'no_transaction_known' description. These attributes and its key codes are then extracted into a dictionary named missing_keys_dict. Additionally, some column attributes in the azdias dataset, not described in the 'DIAS_Attributes_Values' file but having similar attribute properties to other columns in the azdias dataset are engineered to have similar missing value key codes. Next, the missing values are converted to numpy nan.

The second type of missing value is addressed next. These missing values are defined as the initial missing value. The azdias dataset is iterated through to return a list of initial missing values for every column.

Finally a data frame showing the initial missing, final missing and percentage of missing values for every column is created. The figure below shows the top 30 columns with the highest percentage of missing values sorted in descending order:

	Attribute	initial_missing	final_missing	%_missing
7	ALTER_KIND4	890016	890016	99.9
349	TITEL_KZ	73499	889061	99.8
6	ALTER_KIND3	885051	885051	99.3
76	D19_TELKO_ONLINE_DATUM	0	883018	99.1
34	D19_BANKEN_OFFLINE_DATUM	0	871535	97.8
5	ALTER_KIND2	861722	861722	96.7
71	D19_TELKO_ANZ_12	0	857990	96.3
28	D19_BANKEN_ANZ_12	0	831734	93.3
72	D19_TELKO_ANZ_24	0	826208	92.7
87	D19_VERSI_ANZ_12	0	821289	92.2
75	D19_TELKO_OFFLINE_DATUM	0	819114	91.9
4	ALTER_KIND1	810163	810163	90.9
29	D19_BANKEN_ANZ_24	0	794100	89.1
88	D19_VERSI_ANZ_24	0	777037	87.2
35	D19_BANKEN_ONLINE_DATUM	0	726982	81.6
30	D19_BANKEN_DATUM	0	678331	76.1
1	AGER_TYP	0	677503	76.0
73	D19_TELKO_DATUM	0	665798	74.7
100	EXTSEL992	654153	654153	73.4
80	D19_VERSAND_ANZ_12	0	637972	71.6
83	D19_VERSAND_OFFLINE_DATUM	0	634233	71.2
48	D19_GESAMT_ANZ_12	0	584797	65.6
300	KK_KUNDENTYP	584612	584612	65.6
81	D19_VERSAND_ANZ_24	0	563818	63.3
51	D19_GESAMT_OFFLINE_DATUM	0	558558	62.7
49	D19_GESAMT_ANZ_24	0	505303	56.7
84	D19_VERSAND_ONLINE_DATUM	0	494464	55.5
129	KBA05_BAUMAX	133324	476524	53.5
52	D19_GESAMT_ONLINE_DATUM	0	450995	50.6
82	D19_VERSAND_DATUM	0	437886	49.1

Figure 1: Attributes with percentage of missing values



Based on Figure 2, it is noted that only several columns have missing values accounting for more than 30% of their data. A total of 33 columns have more than 30% missing values in their data. These 33 columns/attributes are removed from the data frame.

Besides this other columns/attributes are also checked and removed for the following reasons:

- 1.) identifier that does not contribute to analysis
- 2.) too many zero values
- 3.) too many categories
- 4.) repeat of other attributes
- 5.) attribute is unknown

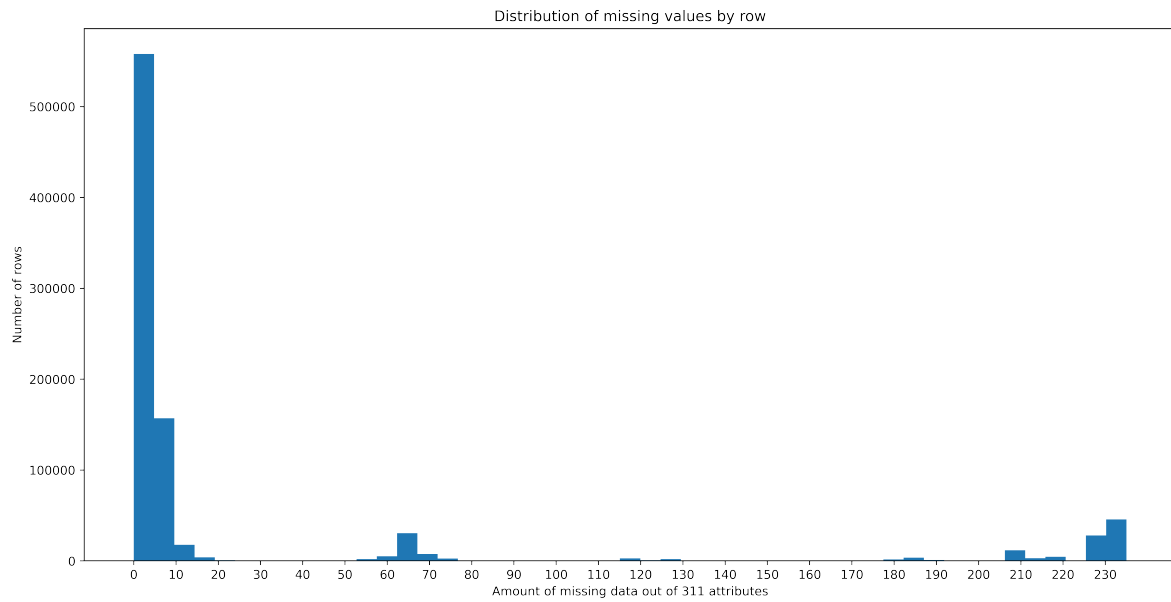


Figure 3: Distribution of missing values by row

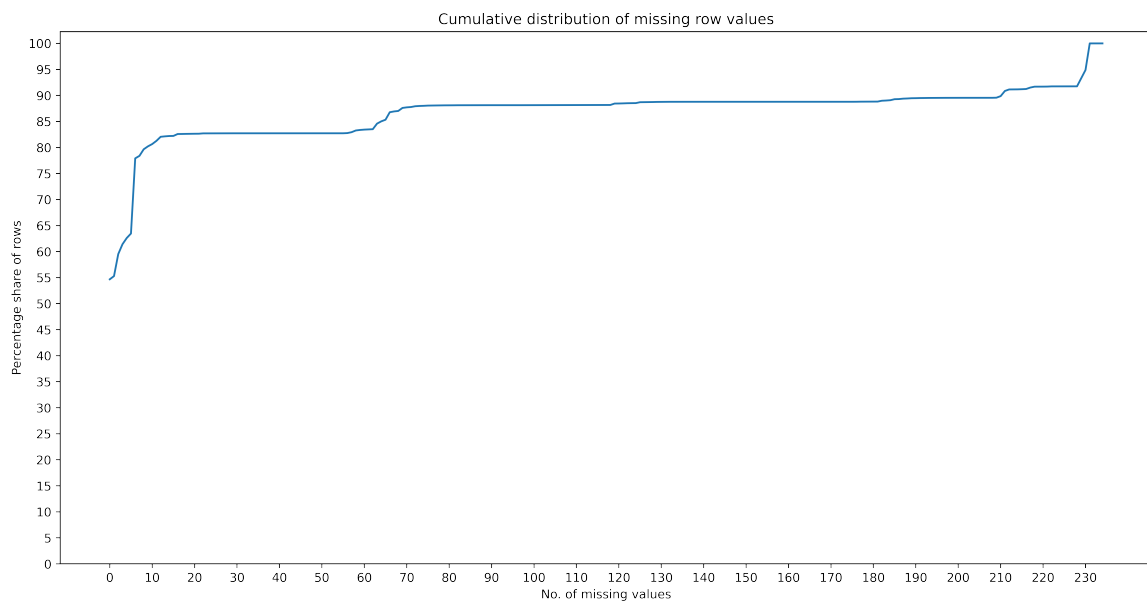


Figure 4: Cumulative distribution of missing row values

The next step is then to assess the number of missing values in rows. In Figure 3, the majority of rows have less than ~20 attributes with missing values.

Based on Figure 4 and 5, up to ~82% of the row data is kept if only rows with less than 26 missing rows are retained. This is the threshold by which row data will be retained.

Upon completing column and row data drops, we perform feature engineering of the azdias and customer dataset.

Feature engineering is then performed on 3 attributes - OST.WEST_KZ, EINGEFUEGT_AM and CAMEO_DEUG_2015.

The remaining rows with missing values is then imputed in the next section.

Percentage of data kept:	
0	54.673869
1	55.290439
2	59.520702
3	61.416080
4	62.605459
5	63.472023
6	77.910978
7	78.394360
8	79.655102
9	80.218711
10	80.656874
11	81.272883
12	82.057649
13	82.139671
14	82.203292
15	82.220684
16	82.588382
17	82.598368
18	82.616209
19	82.629561
20	82.640782
21	82.648299
22	82.711583
23	82.714276
24	82.715398
25	82.718877
26	82.722131
27	82.724711
28	82.727741
29	82.727853

Figure 5: Percentage of row data kept

2.) Impute missing data. Transform and standardise features.

For the remaining missing values, value imputation are performed according to the type of attribute. The imputation performed for the types of attributes are as follows:

- a.) skewed continuous - Attributes that are of continuous nature are checked for its skew value. The pandas 'skew' function is used to identify attributes with skew value of more than 1. A pipeline is then created where these skewed continuous attributes are undergo processing, first being log transformed, then median imputed and finally scaled with standard scaling.
- b.) binary - Attributes that are of binary nature are imputed based on the most-frequent strategy.
- c.) categorical - Attributes that are of categorical nature are imputed based on the most frequent strategy and one-hot encoded.
- d.) numerical - Attributes that are of numerical nature are median imputed and then scaled with standard scaling.

A transformer object is used to combine all the feature transformation pipelines above.

Customer Segmentation

1.) Principal Component Analysis

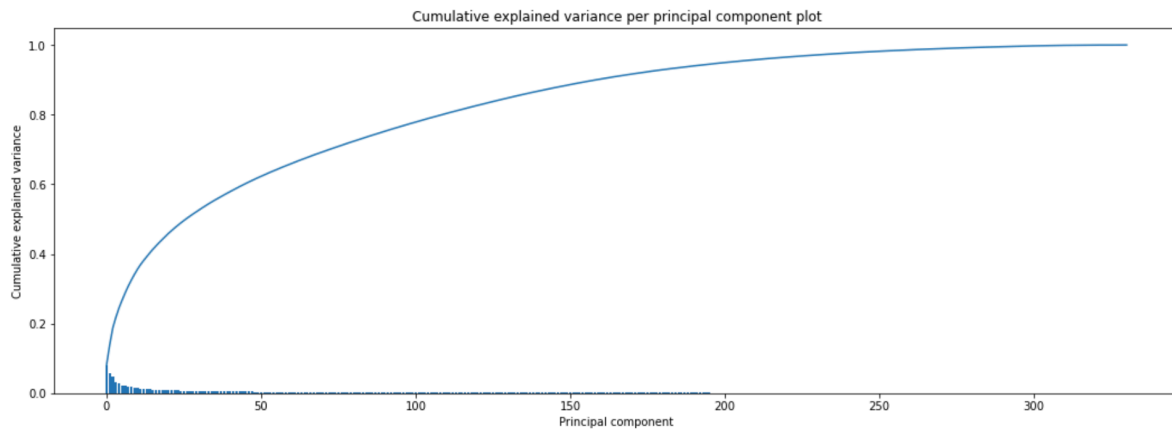


Figure 6: Cumulative explained variance per principal component plot

Principal Component Analysis is used for dimensionality reduction. A plot of cumulative explained variance per principal component plot is produced. From the plot, it is possible to determine that the explained variance value drops off to almost 0 around the 195th component. Here, the cumulative explained variance accounts for 94% of the data variance from the total of 331 features.

By further mapping the PCA components to attributes, we can determine that the top 3 components explains around 18.7% of the data variance.

The first component explains 8.1 % of the data variance and is influenced positively by the attributes ANZ_HAUSHALTE_AKTIV, KBA13_ANTG3, KBA13_ANTG4, KBA13_BAUMAX, PLZ8_BAUMAX. Thus the first component is associated with the number of the household and size of household. The smaller number of household in a building, the more likely they are to be a customer. This could mean single family household in a building is associated more with this component. The larger the size of the household, the more it is associated with this component as well. This could indicate that a single large family household living in a single large building is more associated with this component. A large single family household in a large house means more people in the house that could be a customer and that the customer is likely to be wealthy. The negative attributes influencing this component is MOBI_RASTER, LP_STATUS_GROB, KBA05_ANTG1, KBA13_ANTG1 and MOBI_REGIO. These attributes are inversely associated with this component. Hence, targets with a low social status, high mobility (meaning they likely do not live in an owner-occupier house and is renting, thus frequently moving), and small households (fewer customers) are less likely to be positively associated with this component.

The second component explains 5.7% of the data variance and is associated with KBA13_HERST_BMW_BENZ, KBA13_SEG_OBEREMITTELKLASSE, KBA13_MERCEDES, KBA13_BMW, KBA13_SITZE_4. These attributes relates to expensive and luxury cars owned. Hence, customers with expensive cars are more likely to be a customer. The attributes negatively associated with the component are KBA13_HALTER_25, KBA13_KMH_180, KBA13_SEG_KLEINWAGEN, KBA13_KMH_140_210 and KBA13_SITZE_5. The component is negatively influenced by attributes where the customer does not own a car, owns a low speed vehicle, owns a speed with few seatings and owns a vehicle that is cheap.

The third component explains 4.9% of the data variance and is associated with PRAEGENDE_JUGENDJAHRE, CJT_TYP_1, FINANZ_SPARER, CJT_TYP_2, FINANZ_ANLEGER. These attributes are associated with money saver and investor. Hence, the customer is more likely to be a money saver and investor. The features negatively associated with this component are CJT_TYP_3, CJT_TYP_4, ALTERSKATEGORIE_GROB, CJT_TYP_5 and FINANZ_VORSORGER. These attributes are associated with financial preparedness. Hence, the component is negatively associated with targets who have low financial preparedness.

2.) KMeans Clustering

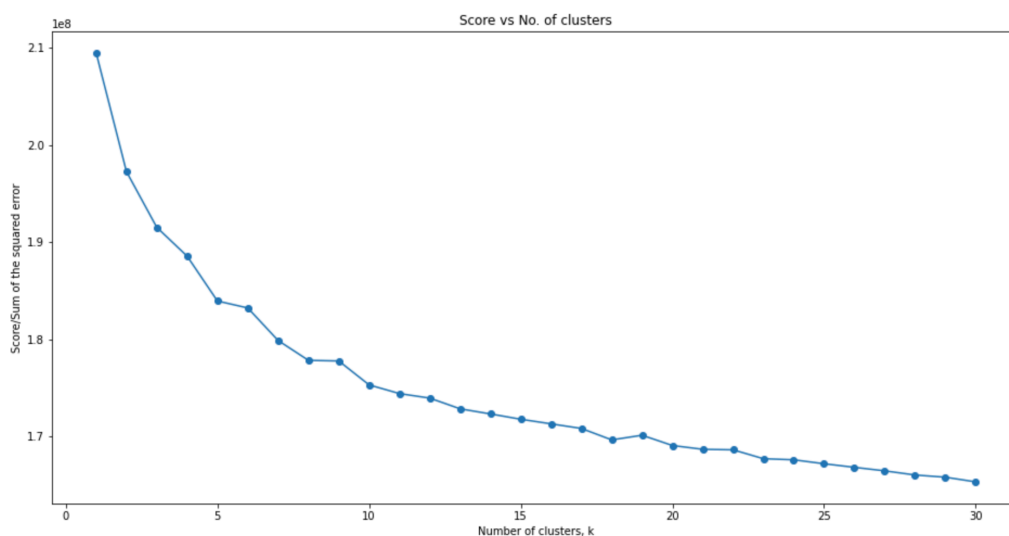


Figure 7: Sum of squared error vs number of clusters

To perform KMeans clustering, an optimal value of 'k' clusters is required. To determine this k-value, an elbow plot will be created. The sum of the squared value roughly plateaus when an optimal number of clusters have been added. This is determined to be at k=8 cluster.

After running the azdias and customers dataset through the clustering pipeline, a plot of percentage share for each cluster is produced for both the azdias (general population) and customers dataset. The resulting plot is as follows:

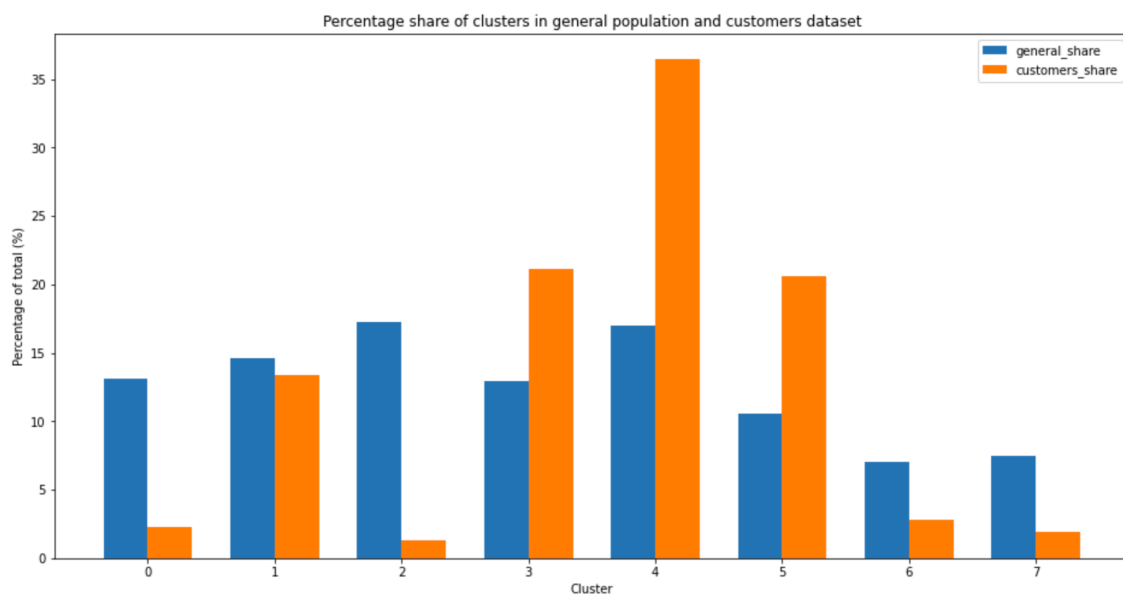


Figure 8: Percentage share of clusters in general population and customers dataset.

From Figure 8, the clusters are represented at a different percentage for the general population and customers dataset. This indicates that a comparison of the percentage representation between the general population and customers dataset can yield insight into the cluster that positively and negatively represents each dataset. Hence, a plot of the percentage difference in total share for general population and customers dataset is created in Figure 9.

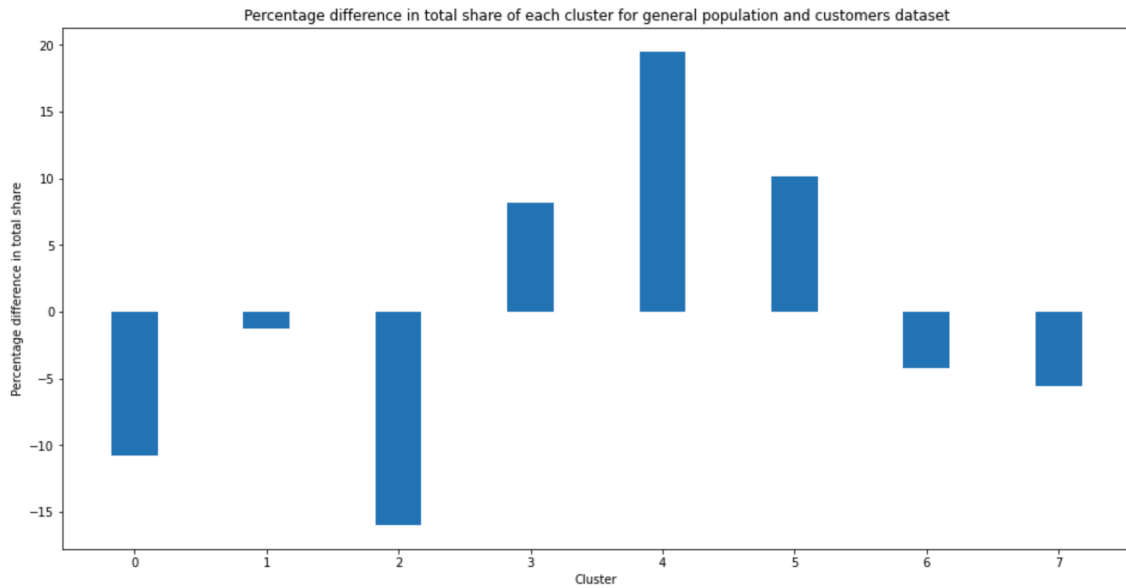


Figure 9: Percentage difference in total share of each cluster between general population and customers dataset

From Figure 9, a positive percentage difference means an over-representation of the customers dataset (compared to the general population dataset), while a negative percentage difference means an under-representation of the customers dataset. The diagram shows prominently that cluster 4 and 5 are over-represented while cluster 0 and 2 are under-represented. The over-represented clusters forms the target group of the marketing campaign and the under-represented clusters is inversely true.

By mapping the clusters back to the attributes of the dataset, an attribute list contributing to the clusters can be studied.

	positive_cluster1	positive_cluster2	negative_cluster1	negative_cluster2
ANZ_STATISTISCHE_HAUSHALTE	2.234183	6.872080	16.062857	5.814821
HH_EINKOMMEN_SCORE	3.428707	3.582609	5.609253	4.862961
SEMIO_REL	2.948076	3.498428	5.197501	5.431387
CJT_TYP_5	4.359608	3.522882	2.178142	2.065598
EWDICHTE	2.862775	4.934159	5.774976	3.588748
FINANZ_ANLEGER	1.957603	1.960719	3.304043	4.076545
FINANZ_SPARER	1.446380	2.240616	3.932492	4.140664
CJT_TYP_2	1.790810	2.866600	4.343586	4.316408

FINANZ_VORSORGER	4.548210	3.795755	2.560997	2.336631
EINGEZOGENAM_HH_JAHR	2000.723413	2003.341237	2006.233312	2006.551229
SEMIO_TRADV	2.778744	3.543177	4.524404	5.101881
KBA13_BAUMAX	1.064956	2.086685	4.201116	1.307635
INNENSTADT	5.608998	3.743530	2.475729	5.137652
MOBI_REGIO	4.125136	3.027193	1.514120	3.012239
GEMEINDE_TYP	31.431368	17.890657	12.577350	26.300017
CJT_TYP_1	2.162840	2.964168	4.258935	4.509262
BALLRAUM	5.050903	2.955694	2.437008	4.569595
ORTSGR_KLS9	3.823973	6.557802	7.781508	4.738989
PRAEGENDE_JUGENDJAHRE	6.138209	8.528607	11.658221	12.730041
CJT_TYP_3	4.366782	3.377998	2.053673	2.275225
AKT_DAT_KL	3.230074	4.202700	5.784368	6.041270
FINANZ_MINIMALIST	4.356588	3.339890	1.813622	2.036623
LP_STATUS_GROB	3.700276	2.792912	1.196847	1.879419
SEMIO_PFLICHT	3.033459	3.712099	5.406316	5.751991
CAMEO_DEUG_2015	3.954690	4.477270	7.980041	5.449069
PLZ8_BAUMAX	1.037183	1.968114	4.046248	1.236788

Table 1: Attributes for over-represented and under-represented clusters

From the table, the attributes of a target customer for the mail order company can be deduced. Attributes of target to be avoided by the company can also be deduced the same way.

Based on inference of the above data, the positively interested clusters (over-represented) are those that have the following attributes:

- financially prepared with savings (FINANZ_SPARER)
- financially lean with good budget (FINANZ_MINIMALIST)
- religiously inclined (SEMIO_REL)
- they are of the 60s - 80s generation (PRAEGENDE_JUGENDJAHRE)
- they have relatively high income (HH_EINKOMMEN_SCORE)
- they are mostly established and consumption minded middle class (FINANZ_MINIMALIST, LP_STATUS, GROB)
- they are slightly above the average income earners (HH_EINKOMMEN_SCORE)

- they are people with investments and are typically money savers (FINANZ_ANLEGER)
- they are traditional-minded (SEMIO_PFLICHT)
- low movement pattern, indicating they live in a house they own (MOBI_REGIO)

The negatively interested clusters (under-represented) are those that have the following attributes:

- very low savings and people who do not invest (FINANZ_ANLEGER)
- low income earners (HH_EINKOMMEN_SCORE, LP_STATUS_GROB)
- less religiously inclined (SEMIO_REL)
- they are of the 80s and above generation and most likely to be a member of the Communist Party Youth Organisation (PRAEGENDE_JUGENDJAHRE)
- they are less traditionally inclined (SEMIO_PFLICHT)
- high movement pattern, indicating the person does not live in a house they own (MOBI_REGIO)

Supervised Learning Model

Through unsupervised learning, the attributes of a most likely customer has been identified. Now supervised learning is applied by building a prediction model that predicts whether a target will be a customer. The MAILOUT_TRAIN dataset will be trained on the model.

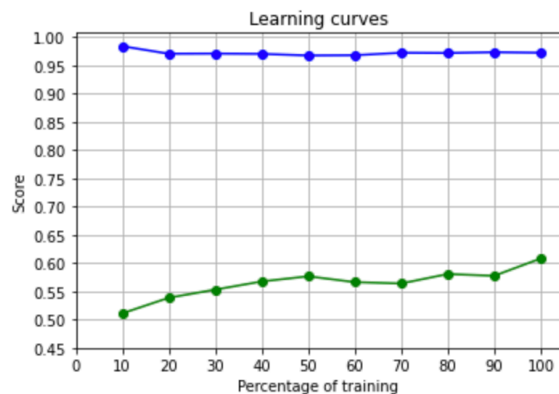
With the supervised learning, the pipeline consists of a column transformer and a classifier. The classifiers chosen are the RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier and XGBClassifier. Ensemble classifiers build on decision-tree model such as RandomForest, AdaBoost and GradientBoosting often improves predictive accuracy by reducing the variation error of unstable classifiers. These classifiers are run on default parameters and the results are then compared to obtain the best classifier.

To evaluate the classifier, we plot the learning curve of the model based on the 'ROC_AUC' score as more data is added to training. This would enable us to determine the classifier that would improve in performance as more data is added and whether which classifier has a high bias or high variance problem.

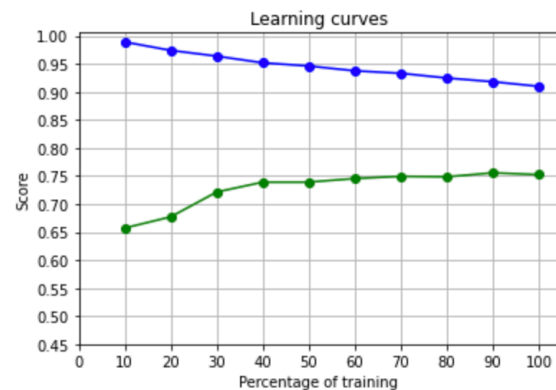
The plots of the learning curve for the classifiers are as follows:

— Training score
— Cross-validation score

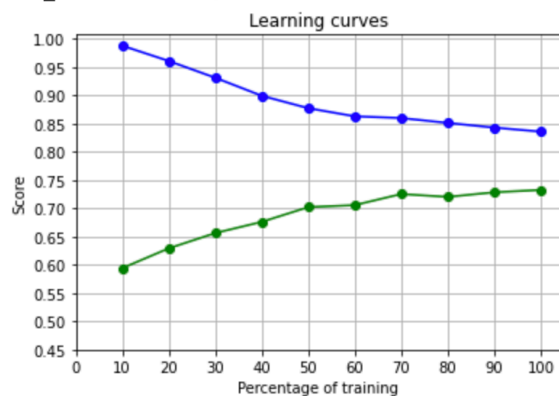
RandomForestClassifier
ROC_AUC Train Score: 0.97
ROC_AUC Validation/Test Score: 0.61



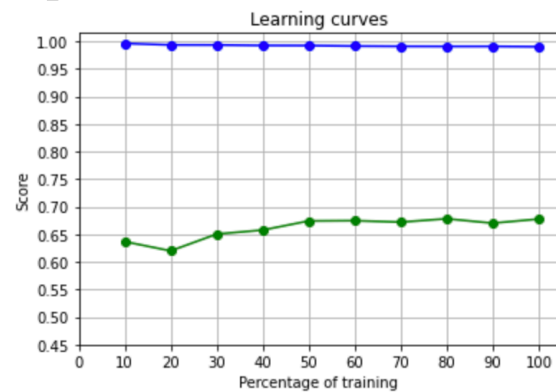
GradientBoostingClassifier
ROC_AUC Train Score: 0.91
ROC_AUC Validation/Test Score: 0.75



AdaBoostClassifier
ROC_AUC Train Score: 0.84
ROC_AUC Validation/Test Score: 0.73



XGBClassifier
ROC_AUC Train Score: 0.99
ROC_AUC Validation/Test Score: 0.68



From the plots, we can see that the RandomForest Classifier and XGB Classifier are a high bias classifier as they both overfit the data.

The AdaBoost Classifier and GradientBoosting Classifier does slightly better as the training and validation curve converges with the addition of data. Both classifiers shows signs that they will converge even more with additional training data. In this case, we select the GradientBoosting Classifier as it has the best validation score.

The GradientBoosting classifier is then parameterised with Grid Search to determine the optimal parameters that would return the best validation score. The optimal parameters for the GradientBoosting classifier is learning_rate: 0.1, max_depth: 3, min_samples_split: 4, n_estimators: 100, with an ROC score of 0.8981.

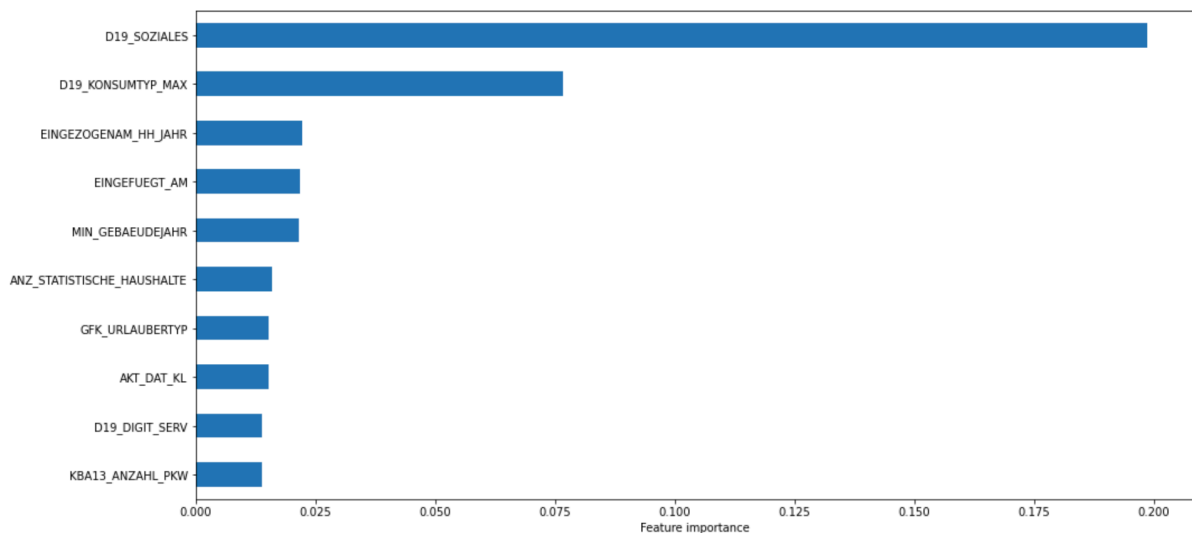


Figure 10: Feature importance plot

Running the estimator to return the ‘feature_importances_’ attribute, it is determined that the estimator relies largely on the attribute of ‘D19_SOZIALES’ and ‘D19_KONSUMTYP_MAX’. The first attribute is not explained in data provided while the second attribute indicates consumption type. However, with an educated guess, it is possible that the first attribute indicates social type.

Kaggle Competition

Finally, the supervised model is used to predict the response of customers based on the MAILOUT_TEST dataset. The resulting score is submitted to the Kaggle competition. The score obtained was 0.78556.

Conclusion

To summarise, the dataset was first pre-processed to handle missing values and a host of other problems, such as attributes not having a description in the data provided. The dataset was then imputed, transformed and scaled with the Column Transformer pipeline. Then unsupervised learning was performed by principal component analysis to reduce the dimensionality of the dataset. Next, Means clustering was used to determine clusters that explain the attributes of a potentially positive target customer. In the supervised learning section, the GradientBoosting classifier was selected as the best classifier. The classifier was then parameterised to obtain the optimal parameters with Grid Search. The

supervised model is then applied to the test data to obtain the final score, which was submitted to Kaggle.

Several improvements can be made to the analysis of this project. Among them includes dropping more columns with missing values, perform feature engineering for more columns, testing different imputation and scaling methods in the Column Transformer pipeline, selecting different classifiers and possibly using over-represented and under-represented cluster attributes datasets only for supervised learning.