

# Machine Learning Engineer Nanodegree

---

## Capstone Proposal

---

Joe Abraham

September 26, 2020

## Proposal

---

*The project is done to fulfill the requirements of Machine learning nanodegree, I am choosing a project and dataset provided by Udacity and Arvato Financial Solutions.*

## Domain Background

### **About Arvato Financial Solutions**

Arvato Financial Solutions offers all financial services related to payments and cash flow. They provide professional financial services to renowned international brands as well as respected local businesses — allowing their customers to leave their credit management to a professional, and they can focus on what matters most for their business.

Their services center around cash flow in all segments of the customer lifecycle: from identity, fraud and credit risk management, to payment and financing services and debt collection.

Source : <https://www.linkedin.com/company/arvato-financial-solutions/about>

### **Project Background**

In this project, Arvato requires a customer segmentation report and wants to predict their potential customers for a mail campaign who will mostly reply back.

## Problem Statement

The Arvato Financial Solutions wants to find the potential customers, additionally they want to know more about the customers and their background. For this reason, analysis of demographics data for customers is done against demographics information for the general population.

1. How to identify parts of the general population to be ideal customers?

2. A method to identify individuals who are likely to be a customer ?

## **Datasets and Inputs**

There are four data files associated with this project:

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).

Additional support files available :

1. DIAS Information Levels - Attributes 2017.xlsx is a top-level list of attributes and descriptions, organized by informational category.
2. DIAS Attributes - Values 2017.xlsx is a detailed mapping of data values for each feature in alphabetical order.

Due to privacy concerns - Datasets are protected and are not publicly available.

## **Solution Statement**

### **Part 0 - Knowing the Data**

Data preprocessing - Inspect the data and necessary data cleaning needs to be done.

### **Part 1 - Unsupervised learning**

Clustering needs to be done with the demographic data provided in Azdias (general population) and customers. For that PCA and K-Means clustering will be the best method and clustering needs to be done with the best model.

### **Part 2 - Supervised Learning**

Here the supervised learning algorithms have to be applied on Udacity\_MAILOUT\_052018\_TRAIN.csv and the best performing model will be tuned for best results.

## Benchmark Model

Supervised Learning, Logistic regression is used as a benchmark model.

## Evaluation Metrics

Performance measurement is an essential task. Since our problem is a classification problem, we can use AUC - ROC Curve to better understand our model. And how our model works. AUC - ROC curve is a performance measurement for classification problems at various thresholds settings. A ROC, or receiver operating characteristic, is a graphic used to plot the true positive rate (TPR, proportion of actual customers that are labeled as so) against the false positive rate (FPR, proportion of non-customers labeled as customers).

Source : <https://www.kaggle.com/c/udacity-arvato-identify-customers/overview/evaluation>

## Project Design

The whole project is divided into 4 parts:

- Part 0 : Know the data, where you try to understand the data
- Part 1 : Customer Segmentation (Unsupervised Learning)
- Part 2 : Supervised Learning
- Part 3 : Kaggle competition

### Part 0 : Know the data, where you try to understand the data

This section contains - Data preprocessing.

- Replacement of strings to float
- Replacement missing values to np.nan - data taken from the excel provided.
- Replacement na with mean
- LNR set as index
- Cleaning columns if it doesn't have enough data ( if the dataset is used for clustering)

### Part 1 : Customer Segmentation (Unsupervised Learning)

This section contains -

- PCA (Dimension Reduction)
- Scaling
- K-Means Clustering with different centroids
- Model fitting and evaluation.

## **Part 2 : Supervised Learning**

- Training in different models and check of their learning curve
- Hyper parameter tuning on the best performing model.
- Model fitting and Scoring.

## **Part 3 : Kaggle competition**

The best performing supervised model prediction is done to the dataset Udacity\_MAILOUT\_052018\_TEST.csv and the results are to be submitted on the kaggle website.