



Home Credit Default Risk

08.14.2018

Rashi Nigam

Springboard Data Science Track

Fremont, CA 94539

Overview

[Home Credit](#) strives to broaden financial inclusion for the unbanked population by providing a positive and safe borrowing experience. In order to make sure this underserved population has a positive loan experience, Home Credit makes use of a variety of alternative data--including telco and transactional information--to predict their clients' repayment abilities.

While Home Credit is currently using various statistical and machine learning methods to make these predictions, they're challenging [Kagglers](#) to help them unlock the full potential of their data. Doing so will ensure that clients capable of repayment are not rejected and that loans are given with a principal, maturity, and repayment calendar that will empower their clients to be successful.

Founded in 1997, Home Credit Group is an international consumer finance provider with operations in 10 countries. It focuses on responsible lending primarily to people with little or no credit history.

Why the Need

- Help better the low income group to get loans who struggle due to insufficient or non-existent credit histories
- Provide a positive and safe borrowing experience to broader community
- Whether influence is in direct relation or inverse relation to TARGET prediction
- Help clients realise their dreams and ambitions in a financially responsible way
- Encourage economic development through supporting domestic consumption, thereby improving living standards

Thus Client needed a research analysis to optimize their service & operations that

- Predict customer ability to repay loan as indicated by TARGET variable in the dataset where TARGET = 0 implies customer is able to repay loan while TARGET = 1 customer's difficulty in repaying loan
- Find correlation among independent feature set that includes exhaustive list of 122 independent features
- Which independent features are influencing the TARGET prediction
- Whether influence is in direct relation or inverse relation to TARGET prediction
- Impact of normalized score of External Source data
- Is there pattern in data distribution among correlated feature set that may help in prediction of TARGET

Goals

Predict how capable each applicant is of repaying a loan. Determine the factors or features that influence Home Credit Default Risk most.

The client is [Home Credit](#) and this research to predict Home Credit Default Risk will be useful to them in knowing:

- What features in the dataset influence the Home Credit Default Risk
- When is the Home Credit Default Risk maximum
- Does other independent features in the dataset like Age, Gender, Employment Duration have any impact on the Home Credit Default Risk? If yes, then is it to advantage or adverse.
- Are External Data influences drive Home Credit Default Risk? If yes, do they have similar influence on Home Credit Default Risk distribution?

Data Wrangling

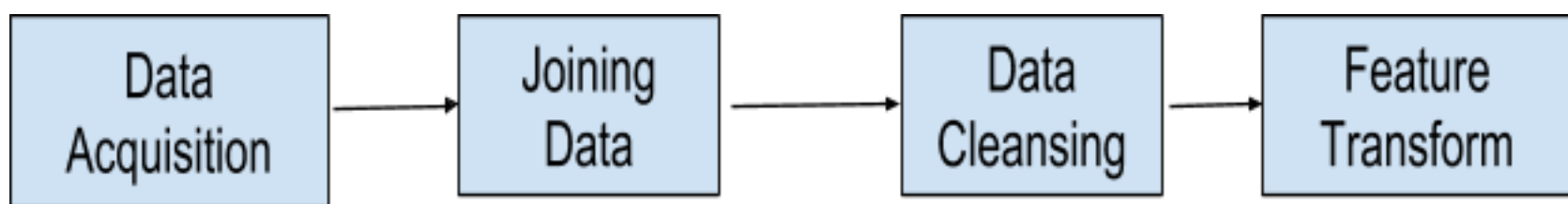
Data wrangling is the process of cleaning and unifying messy and complex data sets for easy access and analysis.

Goals of Data Wrangling are:

- Reveal a “deeper intelligence” within your data, by gathering data from multiple sources

- Provide accurate, actionable data in the hands of business analysts in a timely matter
- Reduce the time spent collecting and organizing unruly data before it can be utilized
- Enable data scientists and analysts to focus on the analysis of data, rather than the wrangling
- Drive better decision-making skills by senior leaders in an organization

Key steps in Data Wrangling



These data are readily and publicly available at

<https://www.kaggle.com/c/home-credit-default-risk/data> and appear as below:

SK_ID_CURR	TARGET	NAME_CONTRACT	CODE_GENDER	FLAG_OWN	FLAG_OWN_CNT	CHILDREN	AMT_INCOME	AMT_CREDIT	AMT_ANNUITY	AMT_GOODWILL	NAME_TYPE	NAME_INCOME	NAME_EDUCATION	NAME_FAMILY	NAME_HOUSE	REGION_POPULATION	DAYS_BIRTH	DAYS_EMPLOYMENT
100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5	351000	Unaccompanied	Working	Secondary / : Single / not r	House / apar	0.018801	-9461	-6	
100003	0	Cash loans	F	N	N	0	270000	1293502.5	35698.5	1129500	Family	State servant	Higher education	Married	House / apar	0.003541	-16765	-11
100004	0	Revolving loans	M	Y	Y	0	67500	135000	6750	135000	Unaccompanied	Working	Secondary / : Single / not r	House / apar	0.010032	-19046	-2	
100006	0	Cash loans	F	N	Y	0	135000	312682.5	29686.5	297000	Unaccompanied	Working	Secondary / : Civil marriage	House / apar	0.008019	-19005	-30	
100007	0	Cash loans	M	N	Y	0	121500	513000	21865.5	513000	Unaccompanied	Working	Secondary / : Single / not r	House / apar	0.028663	-19932	-30	
100008	0	Cash loans	M	N	Y	0	99000	490495.5	27517.5	454500	Spouse, part	State servant	Secondary / : Married	House / apar	0.035792	-16941	-15	
100009	0	Cash loans	F	Y	Y	1	171000	1560726	41301	1395000	Unaccompanied	Commercial	Higher education	Married	House / apar	0.035792	-13778	-31
100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075	1530000	Unaccompanied	State servant	Higher education	Married	House / apar	0.003122	-18850	-4
100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5	913500	Children	Pensioner	Secondary / : Married	House / apar	0.018634	-20099	3652	

The datasets required extensive data wrangling for it involved not only fundamental steps of data preparation but also feature engineering and data imputation to be run during Machine Learning:

1. Extracting Data
2. Identifying Target Dataset among Multiple Data Sources
3. Identifying Missing Data
4. Identifying Data Types of the Feature Set into Non-Categorical and Categorical
5. Casting Data Types per Need
6. Feature Engineering (Date timestamp, One Hot Encoding)
7. Baseline Machine Learning Modeling

Here is the detailed codebook showing above steps:

<https://github.com/rashi-n/Machine-Learning-Projects/blob/master/Capstone%20Projects/Capstone%20II%20Project/Data%20Wrangling.ipynb>

Following the highlights from the codebook:

- I. Among multiple datasets for Credit Accounts and Transactions, 'Application_train' and 'Application_Test' datasets were deduced to import as they seemed to carry most of the features that may be significant in analysis & predictions of the project objective. The main training and testing data with information about each loan application at Home Credit. Every loan has its own row and is identified by the feature SK_ID_CURR. The training application data comes with the TARGET indicating 0: the loan was repaid or 1: the loan was not repaid.:

```
In [2]: #Importing the dataset
df_train = pd.read_csv('application_train.csv')
df_test=pd.read_csv('application_test.csv')
```

```
In [3]: #Shape of dataset
df_train.shape
```

```
Out[3]: (307511, 122)
```

- II. Dataset carries nearly 307K records with 121 features and one 'TARGET' variable to infer predictions on each transaction.
- III. There were 67 columns that had missing values

Your selected dataframe has 122 columns.
There are 67 columns that have missing values.

	Missing Values	% of Total Values
COMMONAREA_MEDI	214865	69.9
COMMONAREA_AVG	214865	69.9
COMMONAREA_MODE	214865	69.9
NONLIVINGAPARTMENTS_MEDI	213514	69.4
NONLIVINGAPARTMENTS_MODE	213514	69.4
NONLIVINGAPARTMENTS_AVG	213514	69.4
FONDKAPREMONT_MODE	210295	68.4
LIVINGAPARTMENTS_MODE	210199	68.4
LIVINGAPARTMENTS_MEDI	210199	68.4
LIVINGAPARTMENTS_AVG	210199	68.4
FLOORSMIN_MODE	208642	67.8
FLOORSMIN_MEDI	208642	67.8
FLOORSMIN_AVG	208642	67.8
YEARS_BUILD_MODE	204488	66.5
YEARS_BUILD_MEDI	204488	66.5
YEARS_BUILD_AVG	204488	66.5
OWN_CAR_AGE	202929	66.0
LANDAREA_AVG	182590	59.4
LANDAREA_MEDI	182590	59.4
LANDAREA_MODE	182590	59.4

IV. Of the 122 features, 106 were non-categorical and 16 were categorical:

```
# Number of each type of column  
df_train.dtypes.value_counts()
```

```
float64      65  
int64        41  
object       16  
dtype: int64
```

V. Feature Engineering

The features were appropriately cast into right data types and categorical features were Label Encoded and One Hot Encoded. The categorical variables with less number of classes (≤ 2) may be label encoded while remaining I may use One Hot encoding technique. The reason why we have to feature engineer Categorical variables is because many Machine Learning Model do not perform well with Categorical variables since categorical variable may have too many levels and pulls down performance level of the model. Also it is not necessary that all the levels always occur in the records, only few may occur and make any impact on model fit.

VI. Aligning Training and Testing Data

There need to be the same features (columns) in both the training and testing data. After feature engineering categorical variables in the dataset now we have more columns in the training (primary) dataset and test dataset since does not have target label so at this point I align the feature set in the two datasets before we take further steps on baselining model.

```
('Training Features shape: ', (307511, 240))
```

```
('Testing Features shape: ', (48744, 239))
```


The training and testing datasets now have the same features which is required for machine learning. The number of features has grown significantly due to one-hot encoding.

VII. Baseline Machine Learning Modeling

With 70 % data missing in some of the 67 columns (total 122 features) it is important to understand how to deal with missing data. As we learn more data science/statistics, we'll learn about data imputation. Here, we'll learn to find missing data points and then we'll drop those points from the dataset so as not to affect our analysis with bias: an important part of data wrangling and data cleaning. We'll try to find which columns in 'df_train.csv' contain missing values and drop those missing values so you'll have tidy data.

Missing Values Strategy # 1 - Identify Features with Missing Values -> Replace with NaN -> Remove all Features with Missing Value -> Assess Model using Logistic Regression

Missing Values Strategy # 2 - Identify Features with Missing Values -> Replace with NaN -> Impute all Features with Missing Value -> Assess Model using Logistic Regression

Followed both the approaches and created Logistic REgression model with following accuracies respectively:

Accuracy of logistic regression classifier on test set: 0.93 with Strategy #1

Accuracy of logistic regression classifier on test set: 0.92 with Strategy#2

```
y_pred = logreg.predict(X_test)
print('Accuracy of logistic regression classifier on test set: {:.2f}'.format(logreg.score(X_test, y_test)))
Accuracy of logistic regression classifier on test set: 0.92
```

Now our baseline model is ready to start Exploratory & Inferential Data Statistics.

Outlier Analysis

Outlier analysis is done to ensure there is no anomaly in the data that may influence model prediction adversely. So I did outlier analysis on 'Birth_days', 'Days_Employed' to ensure data is within possible range.


```
(df_train['DAYS_BIRTH']/365.0).describe()
```

```
count    307511.000000
mean      43.936973
std       11.956133
min       20.517808
25%       34.008219
50%       43.150685
75%       53.923288
max       69.120548
```

```
Name: DAYS_BIRTH, dtype: float64
```

There is no outlier by the AGE either the low or high end. Same we can confirm for Days of Employment.

Model Workflows

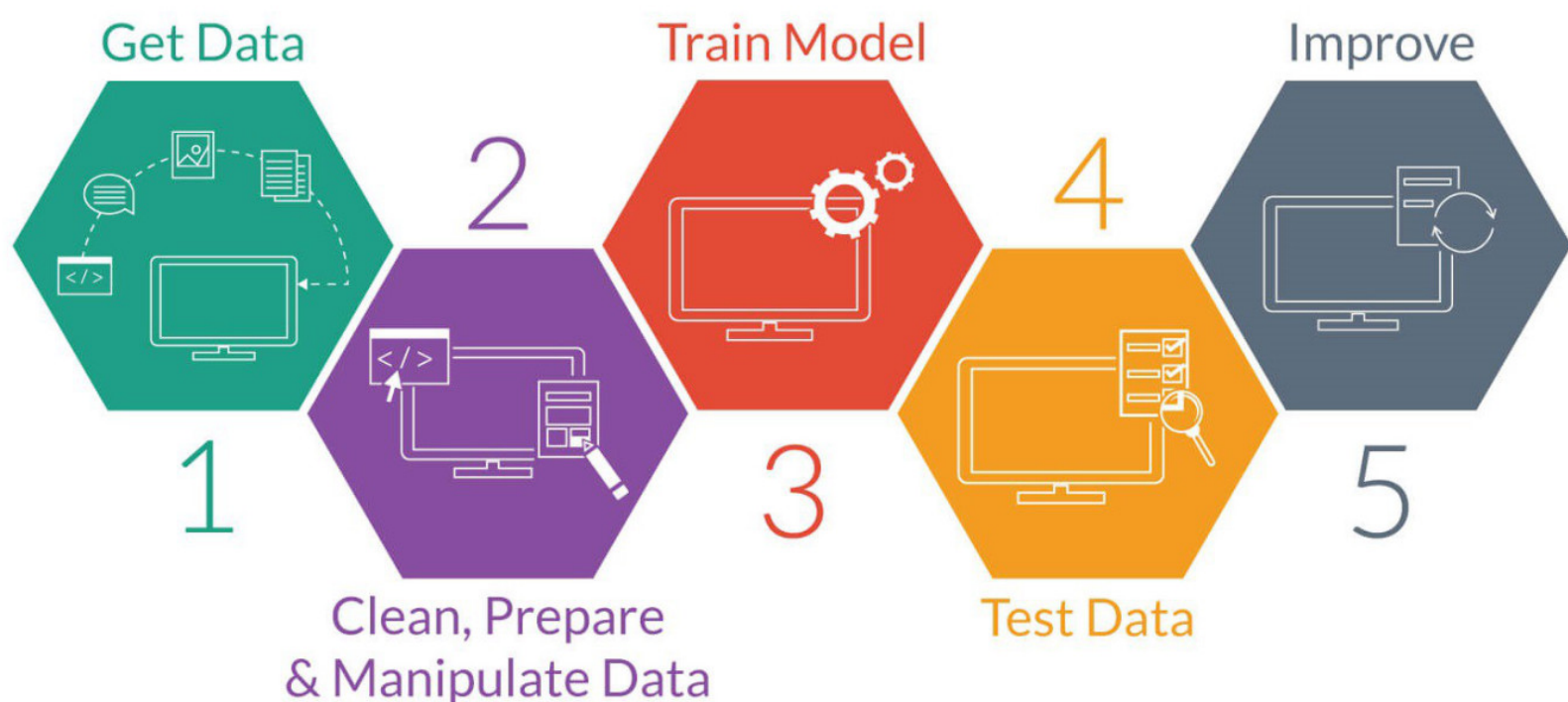


Image courtesy: <https://towardsdatascience.com/building-a-deployable-ml-classifier-in-python-46ba55e1d720>

Any Machine Learning model comprises two stages where Stage I constitutes:

- Get Data/ Data Extraction from authentic/authorized source
- Data Cleaning/Feature Transform/Exploratory Data Analysis
- Train Model using ML techniques after data partitioning into training dataset and Test dataset

Stage II , once training dataset is validated

- Test data is used for the same Machine Learning model
- Compare the results to verify model efficiency and suggest any improvements/recommendations

Exploratory Data Analysis/Inferential Statistics

Of the 122 features, following features showed correlation with the target labeled variable 'TARGET' i.e. 'Home Credit Default Risk':

- Identification if loan is cash or revolving
- Gender of the client
- Normalized score from external data source 1
- Normalized score from external data source 2
- Normalized score from external data source 3
- Flag if the client owns a car
- Flag if client owns a house or flat
- Number of children the client has
- Income of the client
- Credit amount of the loan
- Loan annuity
- For consumer loans it is the price of the goods for which the loan is given
- Who was accompanying client when he was applying for the loan
- Clients income type (businessman, working, maternity leave,...)
- Level of highest education the client achieved
- Family status of the client

- What is the housing situation of the client (renting, living with parents, ...)

Above dependencies can be verified with following visualization graph plots and inferential statistics (code book attached too)

<https://github.com/rashi-n/Machine-Learning-Projects/blob/master/Capstone%20Projects/Capstone%20II%20Project/CSII%20In-Depth%20Analysis.ipynb>

Dataset carries nearly 307K records with 121 features and one 'TARGET' variable to infer predictions on each transaction.

There were 67 columns that had missing values

Of the 122 features, 106 were non-categorical and 16 were categorical:

```
# Number of each type of column
df_train.dtypes.value_counts()

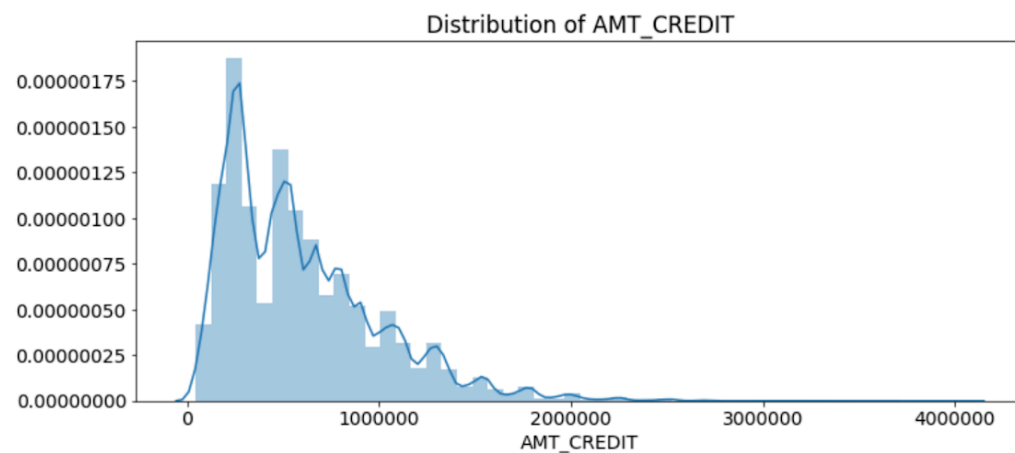
float64    65
int64      41
object     16
dtype: int64
```

1. Exploratory Data Analysis

Visualized distribution of numeric independent features

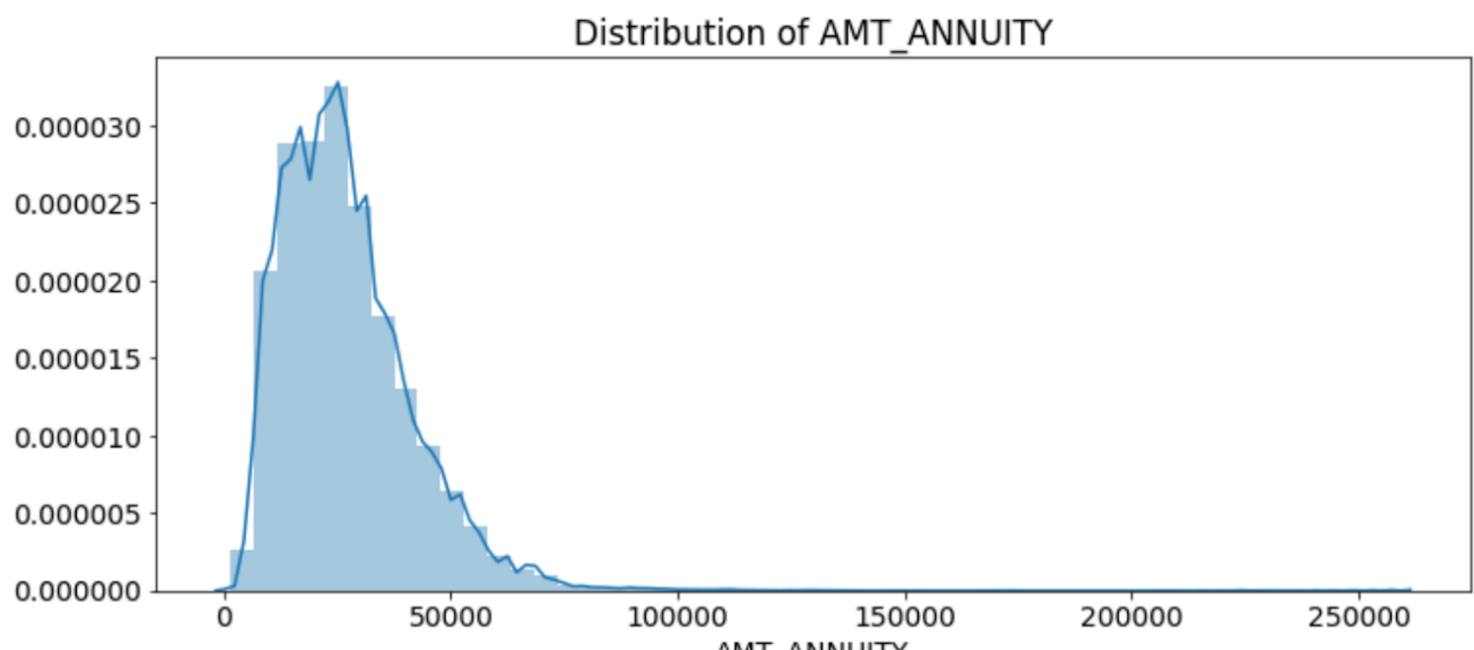
- AMT_CREDIT


```
In [56]: numeric("AMT_CREDIT")
```

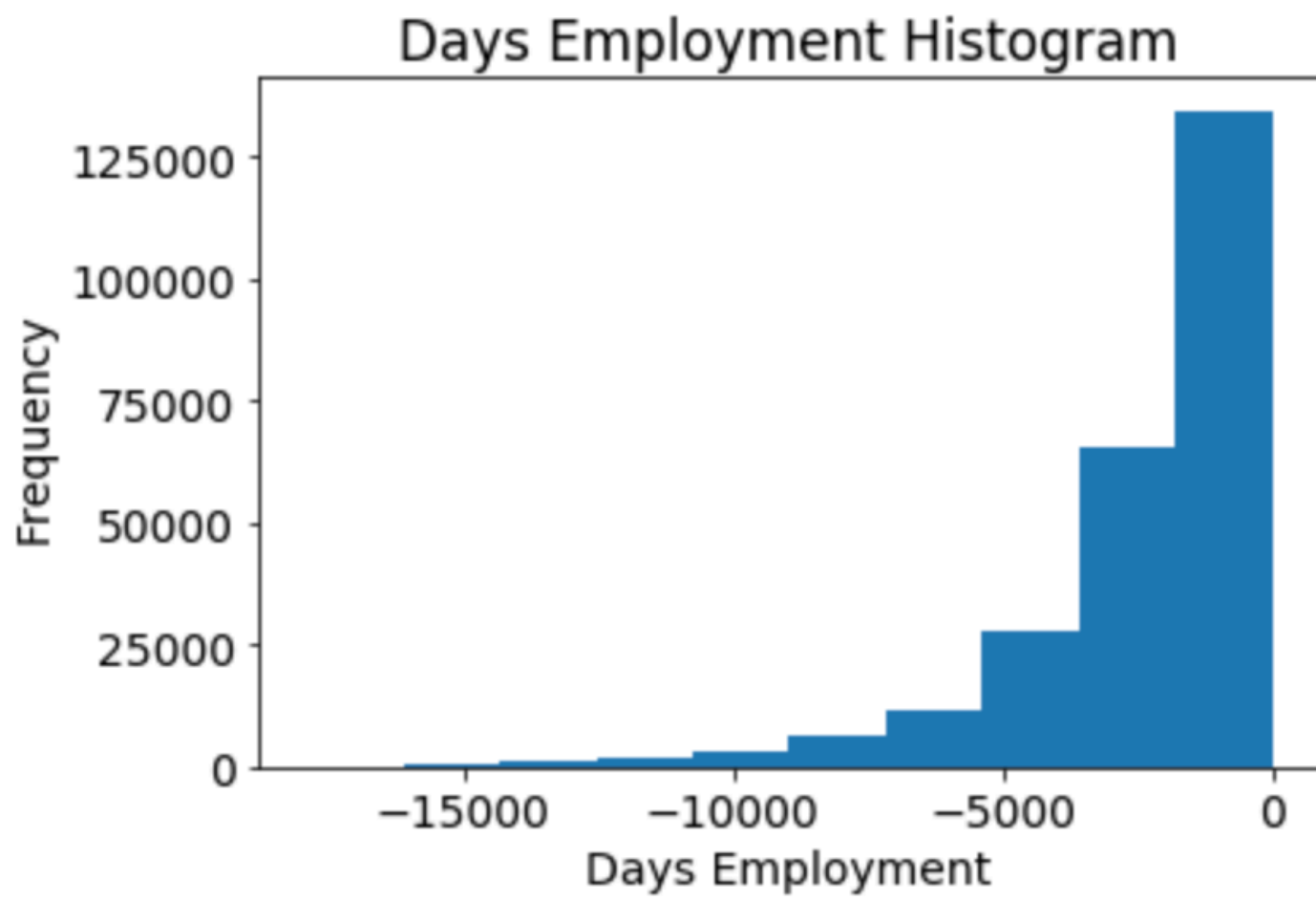


-
- AMT_ANNUIITY

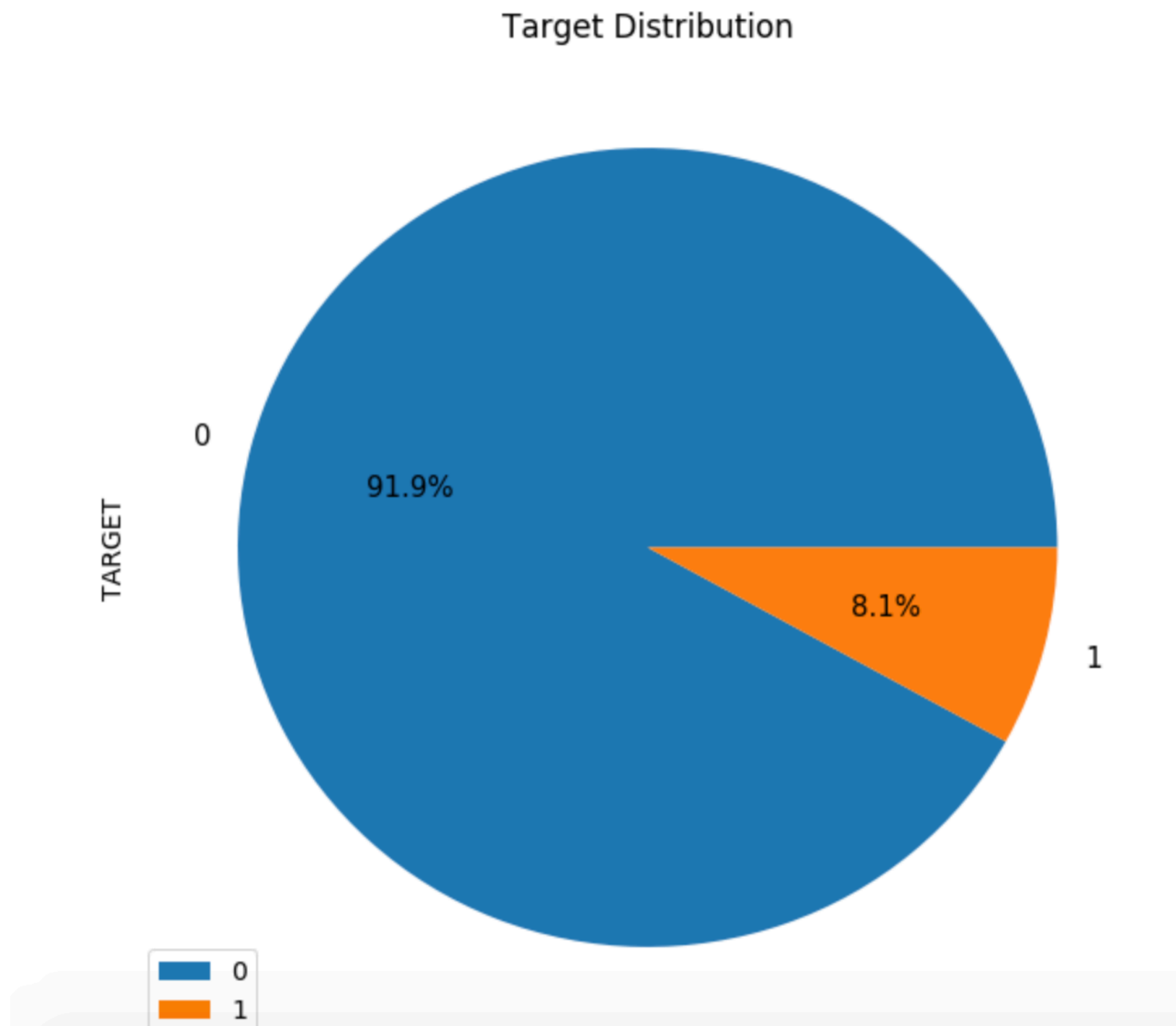
```
numeric("AMT_ANNUIITY")
```



2. The 'Days Employment' distribution looks to be much more in line with what we would expect, and we also have created a new column to tell the model that these values were originally anomalous (because we will have to fill in the nans with some value, probably the median of the column). The other columns with DAYS in the dataframe look to be about what we expect with no obvious outliers. As an extremely important note, anything we do to the training data we also have to do to the testing data. Let's make sure to create the new column and fill in the existing column with np.nan in the testing data. The distribution looks to be much more in line with what we would expect, and we also have created a new column to tell the model that these values were originally anomalous (because we will have to fill in the nans with some value, probably the median of the column). The other columns with DAYS in the dataframe look to be about what we expect with no obvious outliers. As an extremely important note, anything we do to the training data we also have to do to the testing data. Let's make sure to create the new column and fill in the existing column with np.nan in the testing data.



3. Imbalance of Data



From this pie chart, we see this is an imbalanced class problem(<http://www.chioka.in/class-imbalance-problem/>). There are far more loans that were repaid on time than loans that were not repaid. Once we get into more sophisticated machine learning models, we can weight the classes by their representation in the data to reflect this imbalance.

4. Effect of Age on Repayment

By itself, the distribution of age does not tell us much other than that there are no outliers as all the ages are reasonable. Below graph shows more the age of the applicant lesser the risk in client repayment of loan. To visualize the effect of the age on the target, we will next make a kernel density estimation plot (KDE) colored by the value of the target. A kernel density estimate plot shows the distribution of a single variable and can be thought of as a smoothed histogram (it is created by computing a kernel, usually a

Gaussian, at each data point and then averaging all the individual kernels to develop a single smooth curve). We will use the seaborn kdeplot for this graph.

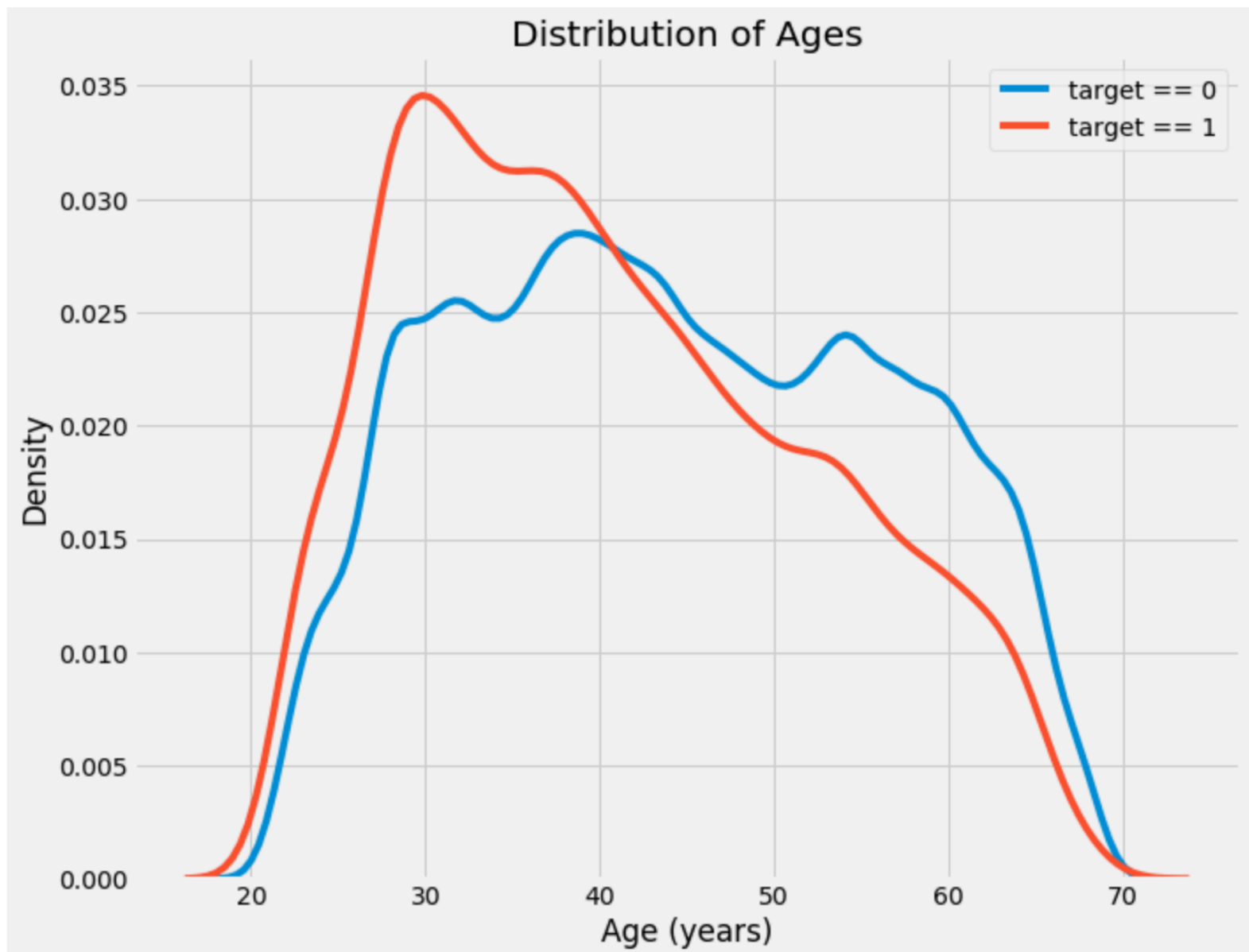
```
# Set the style of plots
plt.style.use('fivethirtyeight')

# Plot the distribution of ages in years
plt.hist(df_train['DAYS_BIRTH'] / 365, edgecolor = 'k', bins = 25)
plt.title('Age of Client'); plt.xlabel('Age (years)'); plt.ylabel('Count');
```



5. The target == 1 curve skews towards the younger end of the range. Although this is not a significant correlation (-0.07 correlation coefficient), this variable is likely going to be useful in a machine learning model because it does affect the target. Let's look at this relationship in another way: average failure to repay loans by age bracket.

To make this graph, first we cut the age category into bins of 5 years each. Then, for each bin, we calculate the average value of the target, which tells us the ratio of loans that were not repaid in each age category.



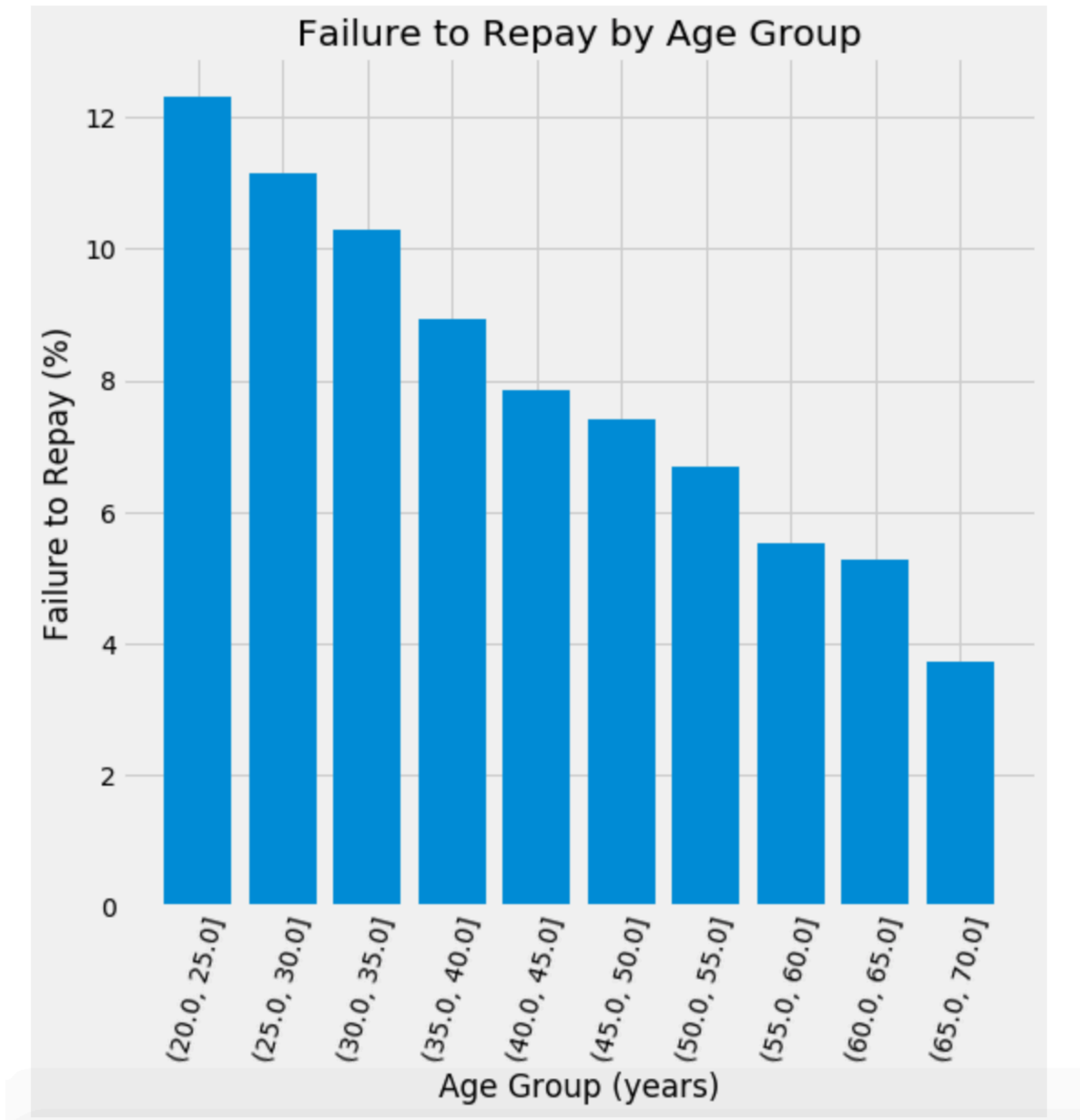
6. Age Group by the Bin & Mean Probability of the TARGET

	TARGET	DAYS_BIRTH	YEARS_BIRTH
YEARS_BINNED			
(20.0, 25.0]	0.123036	8532.795625	23.377522
(25.0, 30.0]	0.111436	10155.219250	27.822518
(30.0, 35.0]	0.102814	11854.848377	32.479037
(35.0, 40.0]	0.089414	13707.908253	37.555913
(40.0, 45.0]	0.078491	15497.661233	42.459346
(45.0, 50.0]	0.074171	17323.900441	47.462741
(50.0, 55.0]	0.066968	19196.494791	52.593136
(55.0, 60.0]	0.055314	20984.262742	57.491131
(60.0, 65.0]	0.052737	22780.547460	62.412459
(65.0, 70.0]	0.037270	24292.614340	66.555108

7. Failure to Repay by Age Group

There is a clear trend: younger applicants are more likely to not repay the loan! The rate of failure to repay is above 10% for the youngest three age groups and below 5% for the oldest age group.

This is information that could be directly used by the bank: because younger clients are less likely to repay the loan, maybe they should be provided with more guidance or financial planning tips. This does not mean the bank should discriminate against younger clients, but it would be smart to take precautionary measures to help younger clients pay on time.

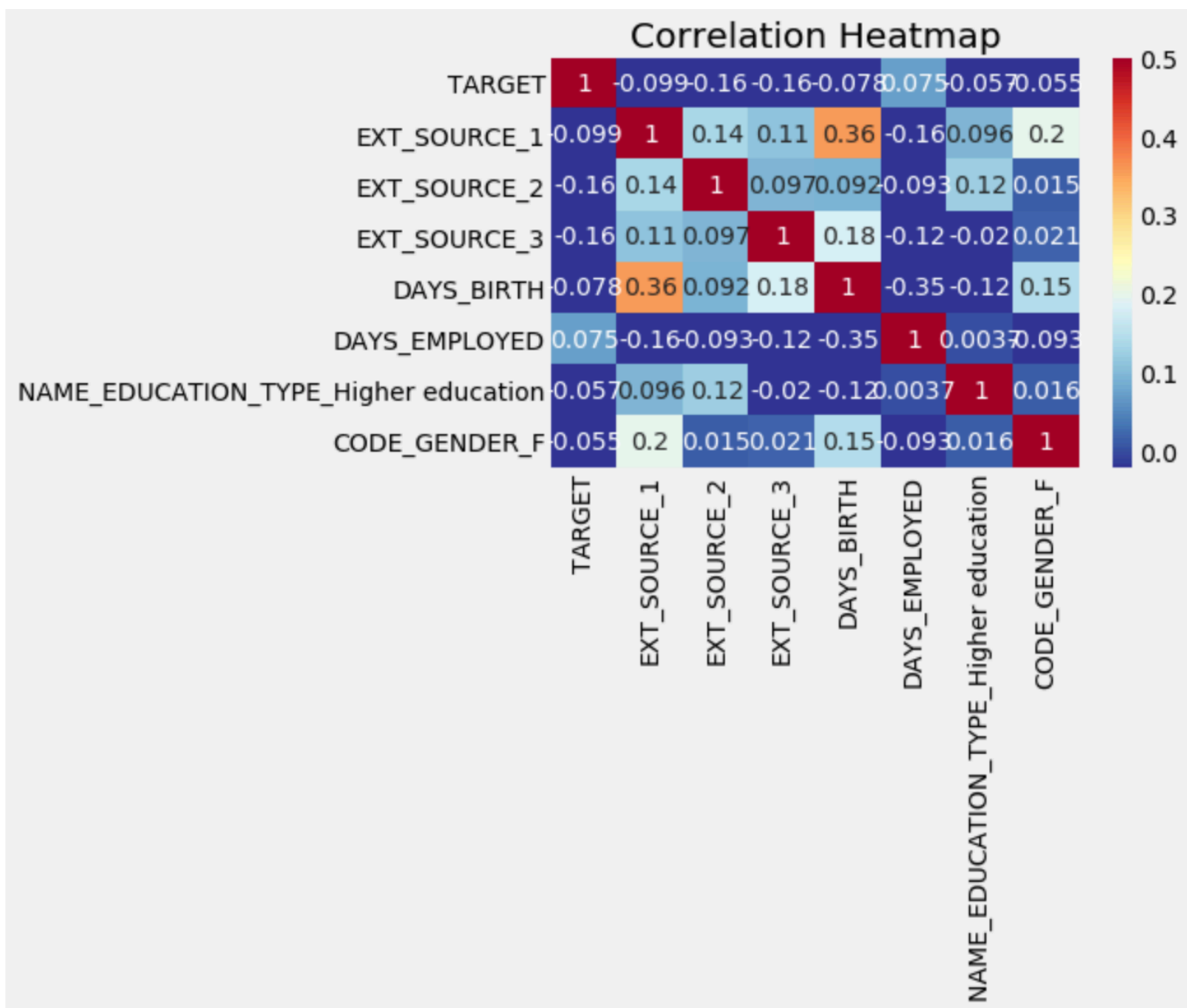


8. Exterior Sources

Exterior Sources The 3 variables with the strongest negative correlations with the target are EXT_SOURCE_1, EXT_SOURCE_2, and EXT_SOURCE_3. According to the documentation, these features represent a "normalized score from external data source". I'm not sure what this exactly means, but it may be a cumulative sort of credit rating made using numerous sources of data.

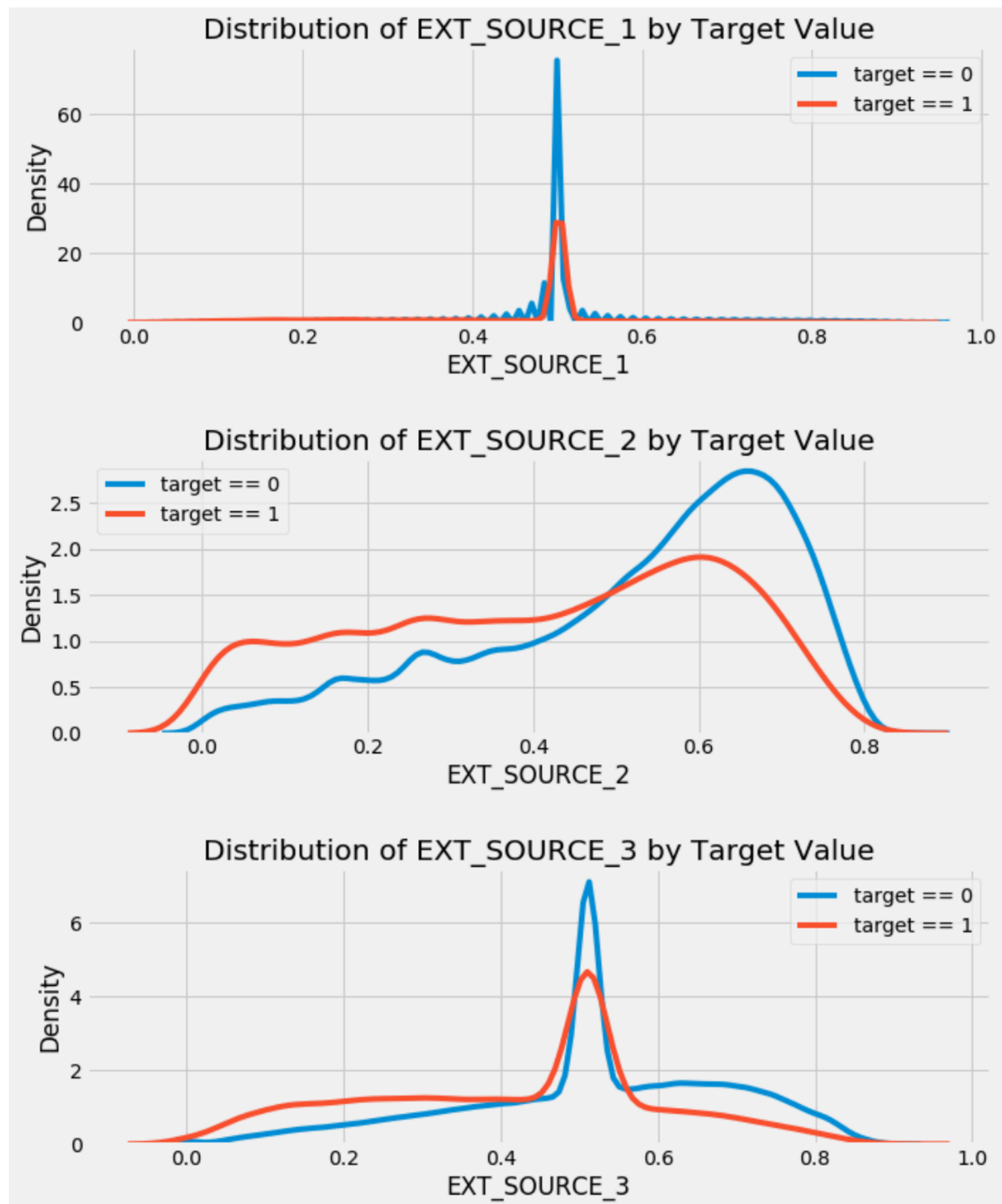
Let's take a look at these variables.

First, we can show the correlations of the EXT_SOURCE features with the target and with each other.



Based on statistical analysis exploring the strength of relationships between TARGET and independent variables such as age, gender, employment duration, external data source, we uncovered the following insights:

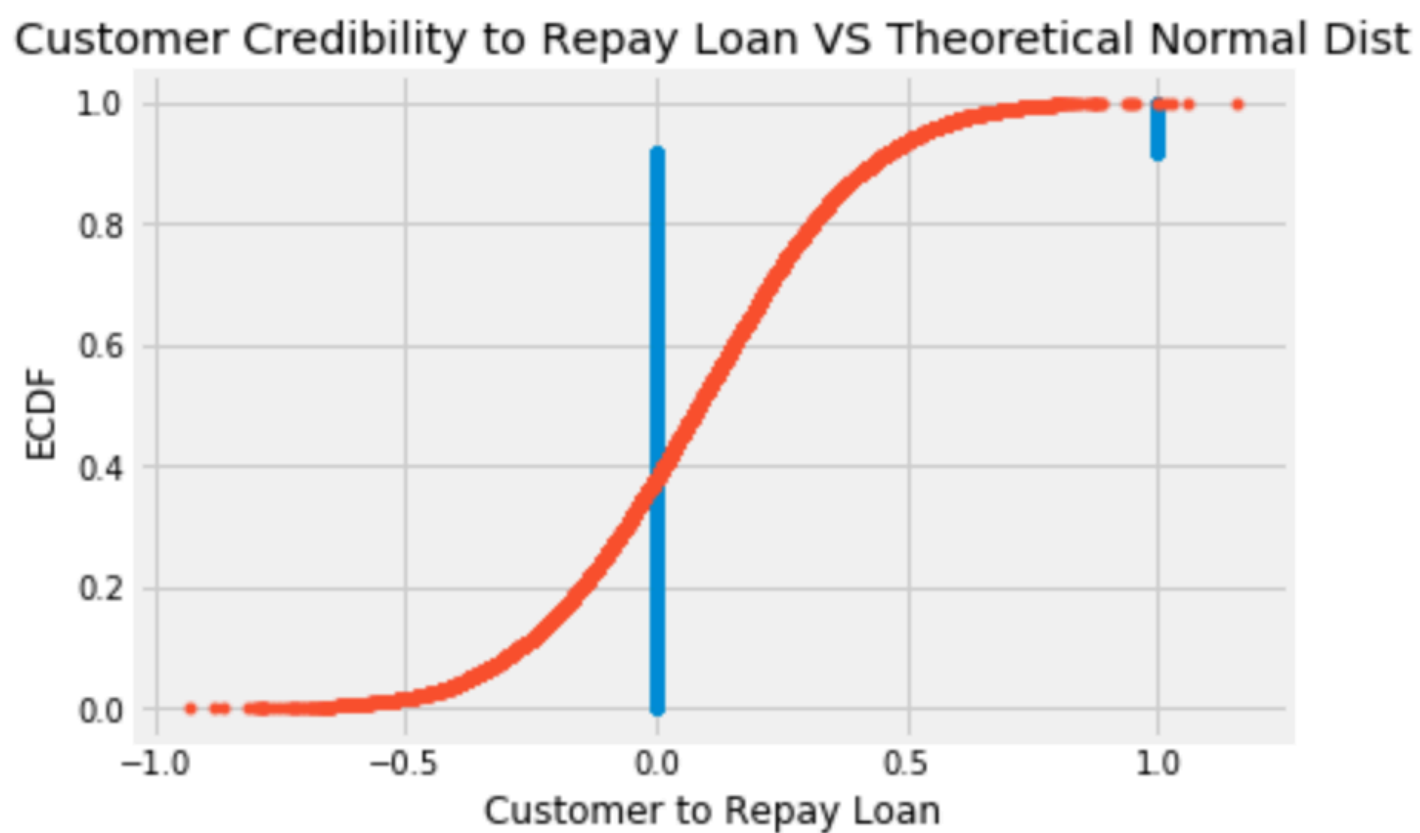
- Age distribution indicates by increasing age TARGET prediction to repay loan increases
- External Data Sources influences TARGET inversely
- Gender correlates with Target prediction
- Employment Duration



Based on these findings, what approach are you going to take? How has your approach changed from what you initially proposed, if applicable?

The next step in our analysis will be to apply inferential statistics: Continuous Distribution, Normal distribution test, Bootstrap sample mean, Confidence Interval to assess Bootstrap replicates, Paired Bootstrap test, Null Hypothesis, Alternate Hypothesis, regression modeling to a majority proportion of the historical data, T-Test, p-value. We will be seeking to determine a best line of fit over the data based on selective application of the variables described earlier. Having chosen the optimal arrangement of our variables, we will test the predictive strength of this model on the remaining portion of our data. This will serve as a secondary check and ensure a minimal amount of model predictions are false positives or negatives. Once this testing phase has validated our model, we can confidently plan to apply the model to future TARGET predictions for the next LOAN transactions in application_train Kaggle datasets.

1. Checking ECDF Distribution of TARGET across the Application_train dataset and theoretical samples of data



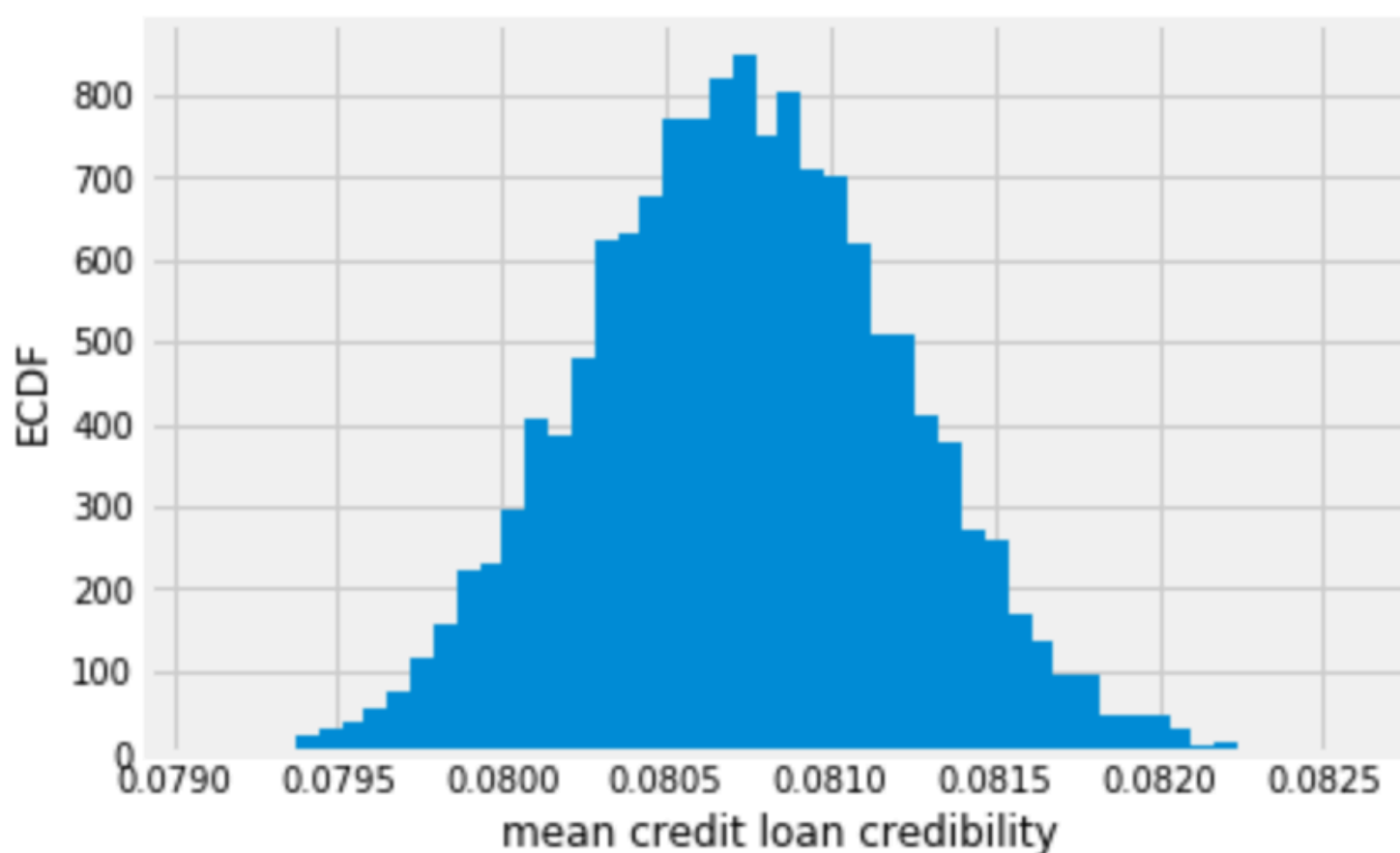
Compare the distribution of the data to the theoretical distribution of the data. This is done by comparing the ecdf First define a function for computing the ecdf from a data set. Next use np.random.normal to sample the theoretical normal distribution and overlay the ecdf of both data sets to compare distribution. Since theoretical ECDF is continuous curve while real data set is contiguous bar for 0 & 1 since it's classification problem but we may consider any data points closer to value '0' indicates 'will repay loan on time', 1 (will have difficulty repaying loan)

2. Confidence Interval

Assuming 95% Confidence interval i.e. give the 2.5th and 97.5th percentile of bootstrap replicates is stored as bs_replicates

```
np.percentile(bs_replicates, [2.5, 97.5])  
O/p: array([0.07979544, 0.08170765])
```

Verifying it with histogram for bootstrap replicates



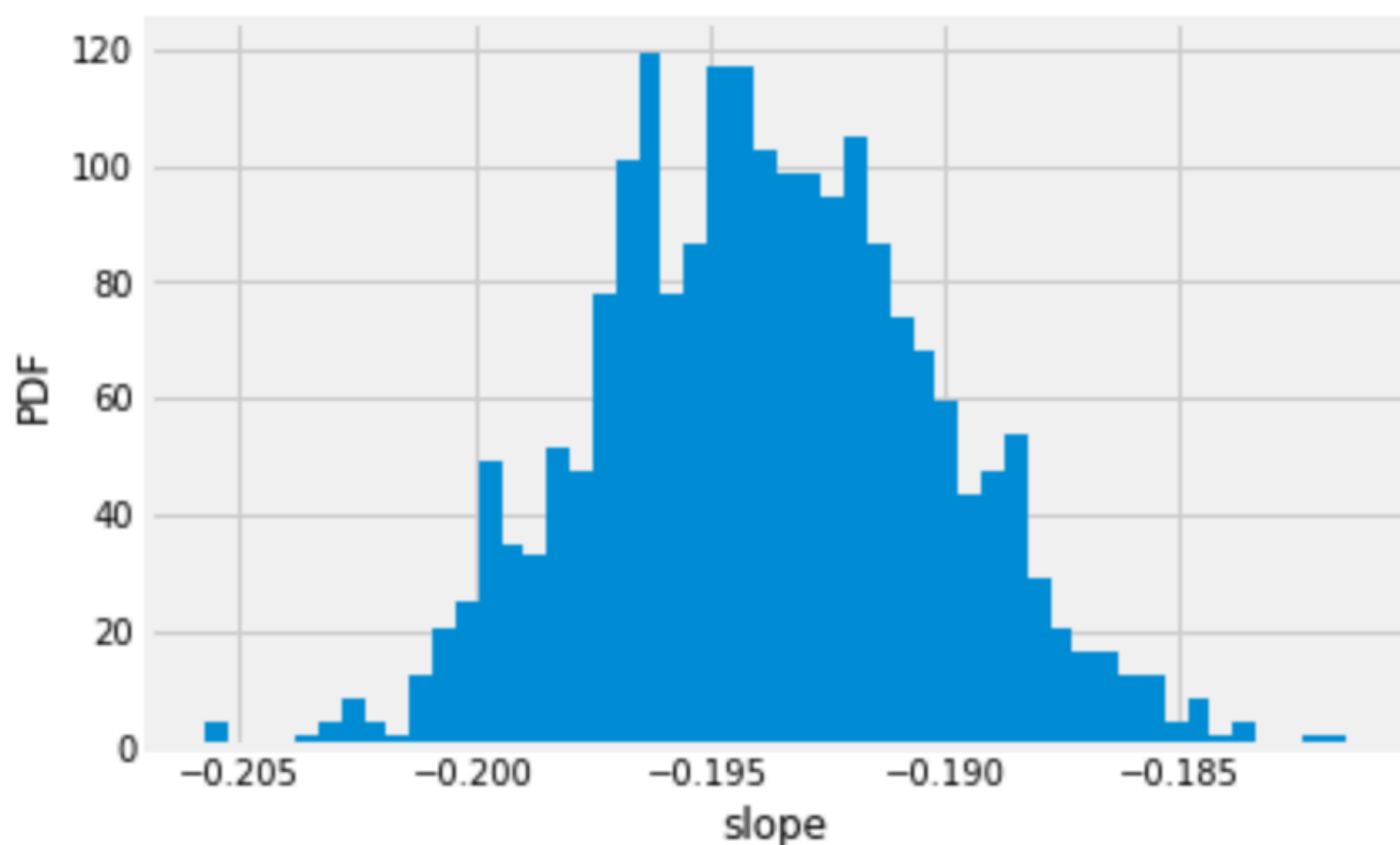
This is bootstrap estimate of the probability distribution function of the mean of 'Credit Loan Default Risk' at the Home Credit Group. Remember, we are estimating the mean 'Credit Loan Default Risk' we would get if the Home Credit Group could repeat all of the measurements over and over again. This is a probabilistic estimate of the mean. I plot the PDF as a histogram, and I see that it is not Normal as it has slightly longer right tail.

In fact, it can be shown theoretically that under not-too-restrictive conditions, the value of the mean will always be Normally distributed. (This does not hold in general, just for the mean and a few other statistics.) The standard deviation of this distribution, called the standard error of the mean, or SEM, is given by the standard deviation of the data divided by the square root of the number of data points. I.e., for a data set.

Notice that the SEM we got from the known expression and the bootstrap replicates is the same and the distribution of the bootstrap replicates of the mean is Normal.

3. Extending Confidence Interval Concept to Pairs Bootstrap

Finding pairs bootstrap for slope & intercept of a function between TARGET (Credit Loan Default risk) and External_Source_1 data



4. Hypothesis Testing

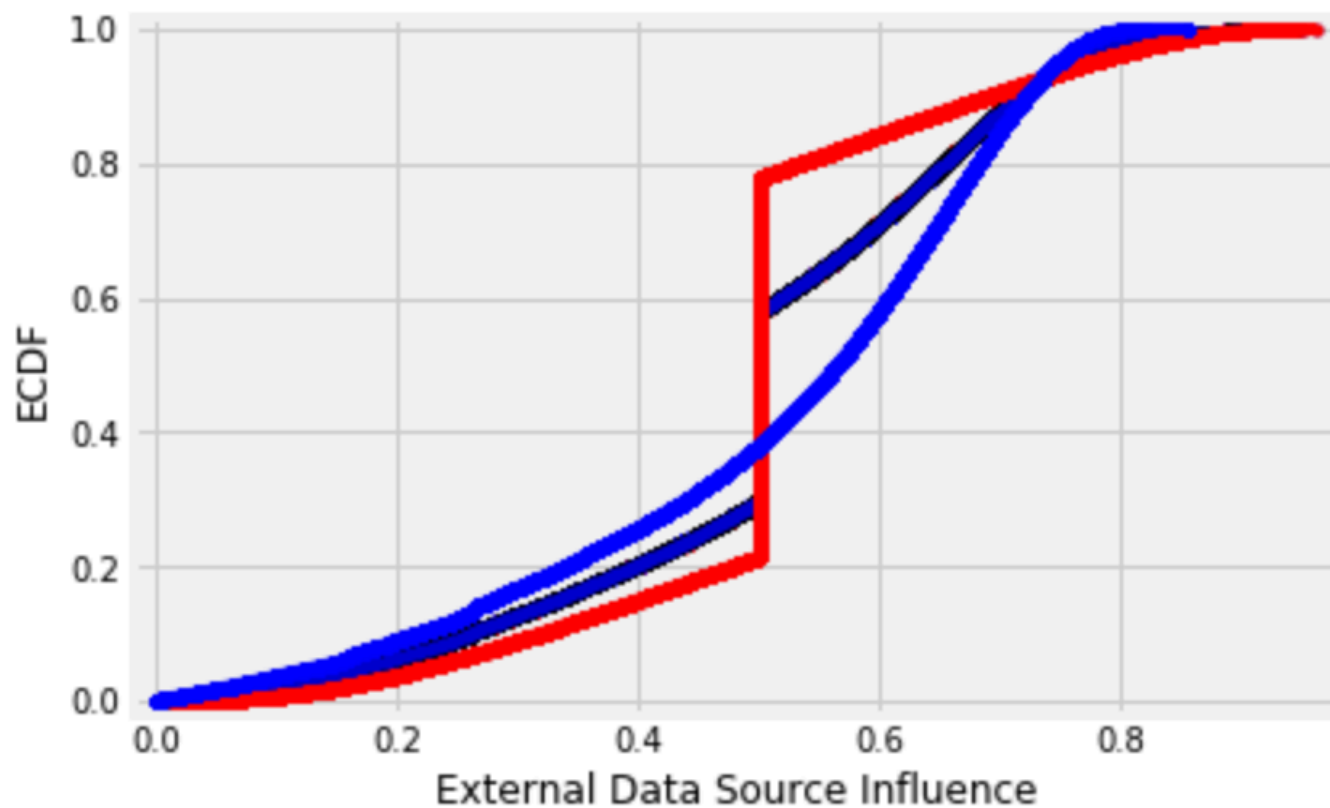
Null Hypothesis- There is no significant difference between EXT_SOURCE_1 and EXT_SOURCE_2 mean on 'Ability to Repay Loan'

$$H_0: \mu_{\text{EXT_SOURCE_1}} - \mu_{\text{EXT_SOURCE_2}} = 0$$

Significance Level: 95% Confidence $\alpha=0.05$

Alternate Hypothesis - There is significant difference between EXT_SOURCE_1 and EXT_SOURCE_2 mean on 'Ability to Repay Loan'

$$H_A: \mu_{\text{EXT_SOURCE_1}} - \mu_{\text{EXT_SOURCE_2}} \neq 0$$



Permutation samples ECDFs overlap and give a purple haze. Few of the ECDFs from the permutation samples overlap with the observed External Source Data1 data towards right of the graph & even fewer overlap towards left, suggesting that the hypothesis is not commensurate with the data. External Source Data1 & External Source Data2 are not identically distributed and do not influence data in similar way. So Null Hypothesis is rejected.

EDA Outcomes

Based on above EDA & Inferential Statistics, following are the outcomes:

- Age distribution indicates by increasing age TARGET probability to repay loan increases
- External Data Sources influences TARGET inversely
- Gender correlates with Target prediction
- Employment Duration indicates by increasing Employment Duration, TARGET probability to repay loan increases

- The feature set that drives this prediction is:
 - Identification if loan is cash or revolving
 - Gender of the client
 - Normalized score from external data source 1
 - Normalized score from external data source 2
 - Normalized score from external data source 3
 - Flag if the client owns a car
 - Flag if client owns a house or flat
 - Number of children the client has
 - Income of the client
 - Credit amount of the loan
 - Loan annuity
 - For consumer loans it is the price of the goods for which the loan is given
 - Who was accompanying client when he was applying for the loan
 - Clients income type (businessman, working, maternity leave,...)
 - Level of highest education the client achieved
 - Family status of the client
 - What is the housing situation of the client (renting, living with parents, ...)

- There is direct positive relation between TARGET (Home Credit Default Risk) vs 'Age Distribution' while negative relation with 'External _Sources'
- TARGET distribution curve is almost similar to Normal distribution curve with slightly longer right tail
- Comparing ECDF Target (Credit Loan Default risk) for the dataset with theoretical samples and bootstrap samples affirms 95% confidence interval.
- Extended Confidence Interval Concept to Pairs Bootstrap using slope & intercept of a function between between TARGET (Credit Loan Default risk) and External_Source_1 data: Null hypothesis i.e. Ext_Source_1 & Ext_Source_2 carry similar influence on Target is rejected.

Here is the link to iPython Notebook:

<https://github.com/rashi-n/Machine-Learning-Projects/blob/master/Capstone%20Projects/Capstone%20II%20Project/EDA%20%26%20Inferential%20Statistics%20II.ipynb>

Machine Learning Modelling

I next propose to do in-depth analysis on the given dataset using one of the Supervised Machine Learning Classification algorithm since output datasets are provided and I can use this to predict the target variable.

Objective to undertake Machine Learning is because it focuses on the design of systems that can learn from and make decisions and predictions based on data. Machine learning enables computers to act and make data-driven decisions rather than being explicitly programmed to carry out a certain task. The iterative aspect of machine learning is important because as models are exposed to new data, they can independently adapt. They learn from previous computations to produce reliable, repeatable decisions and results. It's a science that's not new – but one that has gained fresh momentum.

This project is Classification problem as dependent variable 'TARGET' value represents class to which a data point is part of (discrete value) and can be used to predict the output value using training data.

TARGET is assessment of Home Credit Default Risk and that I analyzed to find correlation, mean, standard deviation, confidence interval, ECDF, permuted paired bootstrap replicate plots to find optimal function for prediction with other related features(labeled data) in the dataset.

Following are few independent variables that are influencing the outcome of TARGET in this project:

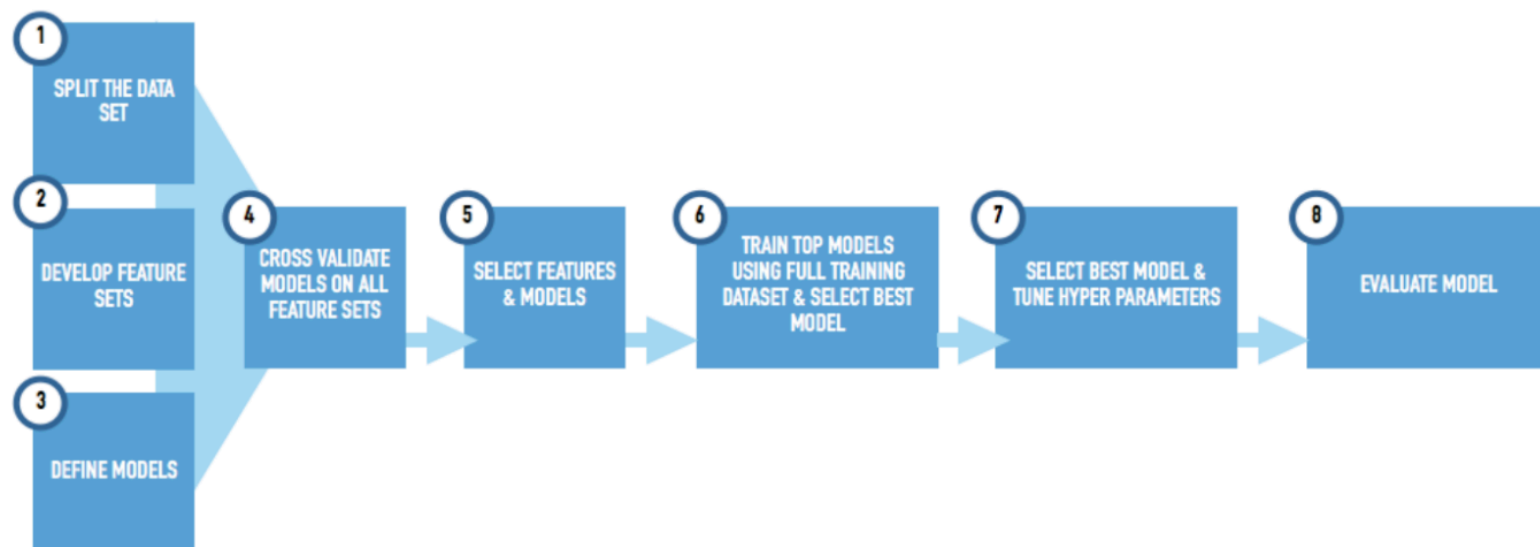
- Age
- External Data Sources Influence Normalized
- Gender
- Employment duration
- Higher Education

Steps taken in Model Analysis are:

Machine learning brings the promise of deriving meaning from all of that data. Machine learning is a problem of induction where general rules are learned from specific observed data from the domain. Machine learning techniques include clustering, where objects are grouped into bins with similar traits; regression, where relationships among variables are estimated; and classification, where a trained model is used to predict a categorical response.

For the purpose of this project, it is Machine Learning Classification problem and following steps are performed while ML Modeling:

- I. Gathering data
- II. Data Preparation (Missing values, Categorical values, Imbalance in Data handled)
- III. Choosing a baseline model
- IV. Training
- V. Evaluation
- VI. Tuning
- VII. Deploy the model for Prediction



Step 1: Gathering Data

This step is very important because the quality and quantity of data that you gather will directly determine how good your predictive model can be. In this case, the data we collect will be the 121 independent features that determine the TARGET variable.

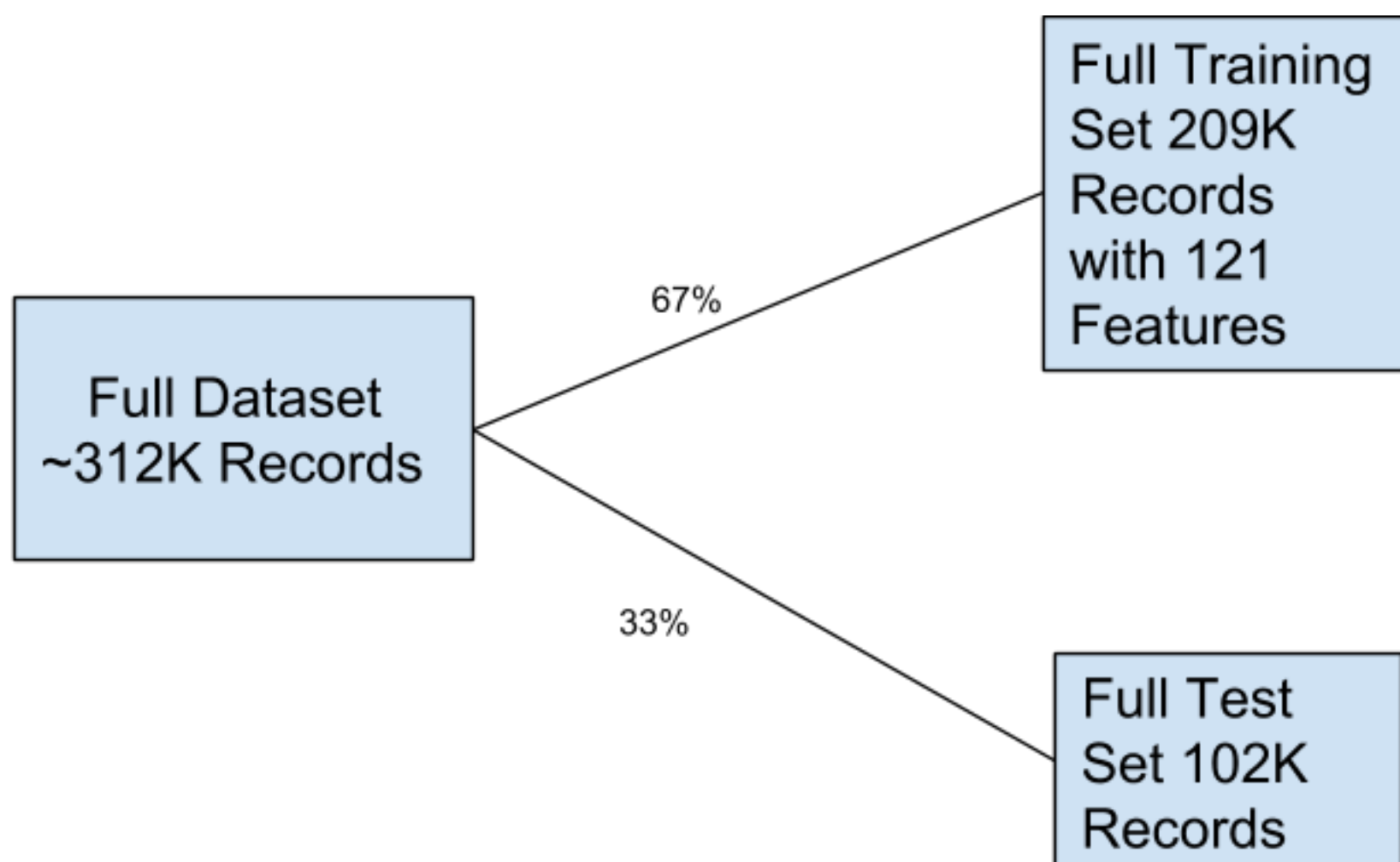
The analysis was done in jupyter notebooks:

<https://github.com/rashi-n/Machine-Learning-Projects/blob/master/Capstone%20Project%20II%20Project/CSII%20In-Depth%20Analysis.ipynb>

Step 2:

Data preparation

Data preparation, where we load our data into a suitable place and prepare it for use in our machine learning training. We'll first put all our data together, and then randomize the ordering. We don't want the order of our data to affect what we learn.



This is also a good time to do any pertinent visualizations of your data, to help you see if there are any relevant relationships between different variables you can take advantage of, as well as show you if there are any data imbalances. For example, if we collected way more data points about Customers with '0' risk to repay loan than not, the model we train will be biased toward guessing that virtually everyone that it sees is able to repay loan, since it would be right most of the time. However, in the real-world, the model may see customers of both classes an equal amount, which would mean that guessing "0" would be wrong half the time.

We'll also need to split the data in two parts. The first part, used in training our model, will be the majority of the dataset. The second part will be used for evaluating our trained model's performance. We don't want to use the same data that the model was

trained on for evaluation, since it could then just memorize the “questions”, just as you wouldn’t use the same questions from your math homework on the exam.

Sometimes the data we collect needs other forms of adjusting and manipulation. Things like de-duping, normalization, error correction, and more. These would all happen at the data preparation step. In this case, I have data preparation needs for missing data, null entries, categorical variables, imbalance in dataset.

Training subset is basically including all the features in predictive analytics.

Step 3:

Choosing Model : The next step in our workflow is choosing a model. There are many models that researchers and data scientists have created over the years. Some are very well suited for image data, others for sequences (like text, or music), some for numerical data, others for text-based data. In our case, since have 121 features, and one TARGET variable we can use a Supervised Classification model, which is a fairly simple one that should get the job done. Three Machine Learning models were analyzed for this Supervised Classification problem:

1. Logistic Regression
2. Random Forest Classification
3. Decision Tree Classification

Step 4:

Training: Now we move onto what is often considered the bulk of machine learning—the training. In this step, we will use our data to incrementally improve our model’s ability to predict whether a given Customer transaction is risk ‘0’ or ‘1’ in terms to repay loan.

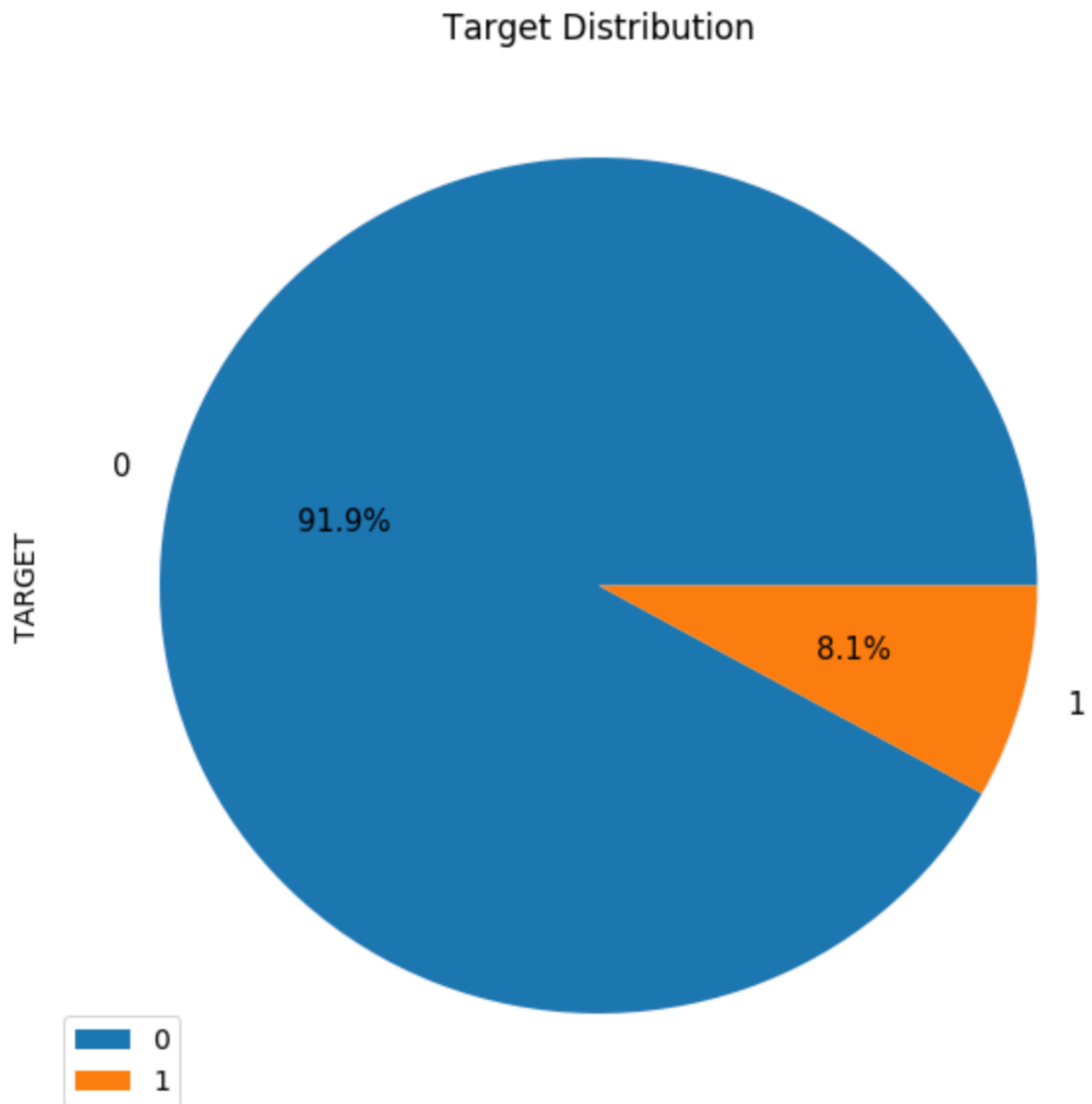
In some ways, this is similar to someone first learning to drive. At first, they don’t know how any of the pedals, knobs, and switches work, or when any of them should be used. However, after lots of practice and correcting for their mistakes, a licensed driver emerges. Moreover, after a year of driving, they’ve become quite adept. The act of driving and reacting to real-world data has adapted their driving abilities, honing their skills.

Logistic Regression/Random Forest/Decision Tree Classification was performed on all feature set with One-Hot/Label Encoding applied for Categorical variables, Missing Data strategies (removing rows with Null values AND Imputation) and Imbalance in dataset(resampling, TOMEK links, SMOTE)

There were 16 Categorical Variables and ones with relatively small number of unique entries were treated with label encoding:

NAME_CONTRACT_TYPE	2
CODE_GENDER	3
FLAG_OWN_CAR	2
FLAG_OWN_REALTY	2
NAME_TYPE_SUITE	7
NAME_INCOME_TYPE	8
NAME_EDUCATION_TYPE	5
NAME_FAMILY_STATUS	6
NAME_HOUSING_TYPE	6
OCCUPATION_TYPE	18
WEEKDAY_APPR_PROCESS_START	7
ORGANIZATION_TYPE	58
FONDKAPREMONT_MODE	4
HOUSETYPE_MODE	3
WALLSMATERIAL_MODE	7
EMERGENCYSTATE_MODE	2

Imbalance in Dataset could be easily visualized using following EDA:



Based on above visualization we can easily infer that this is highly biased data.

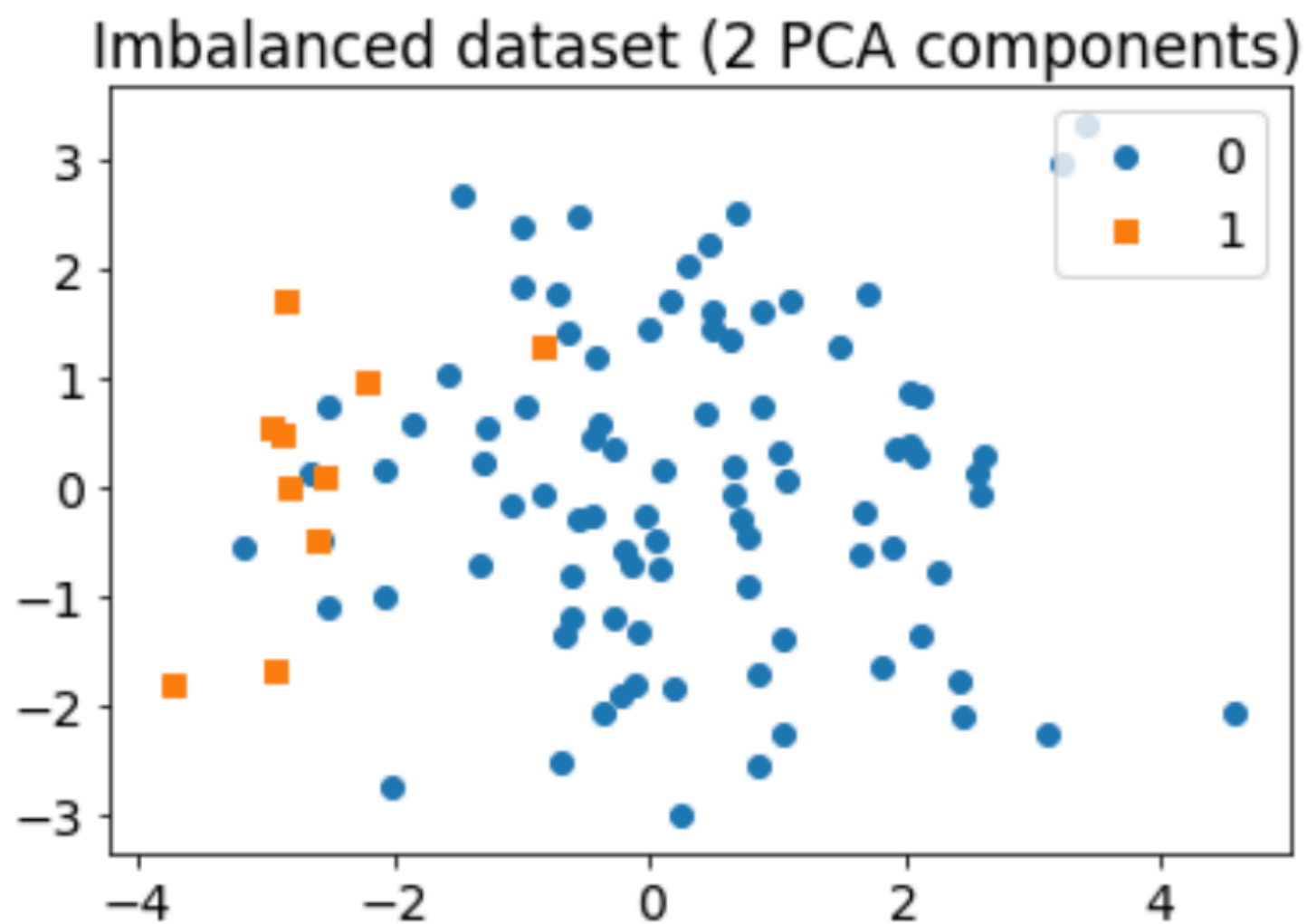
Resampling is the method that consists of drawing repeated samples from the original data samples. The method of Resampling is a nonparametric method of statistical inference. In other words, the method of resampling does not involve the utilization of the generic distribution tables (for example, normal distribution tables) in order to compute approximate p probability values.

Resampling involves the selection of randomized cases with replacement from the original data sample in such a manner that each number of the sample drawn has a number of cases that are similar to the original data sample. Due to replacement, the drawn number of samples that are used by the method of resampling consists of repetitive cases. Courtesy:

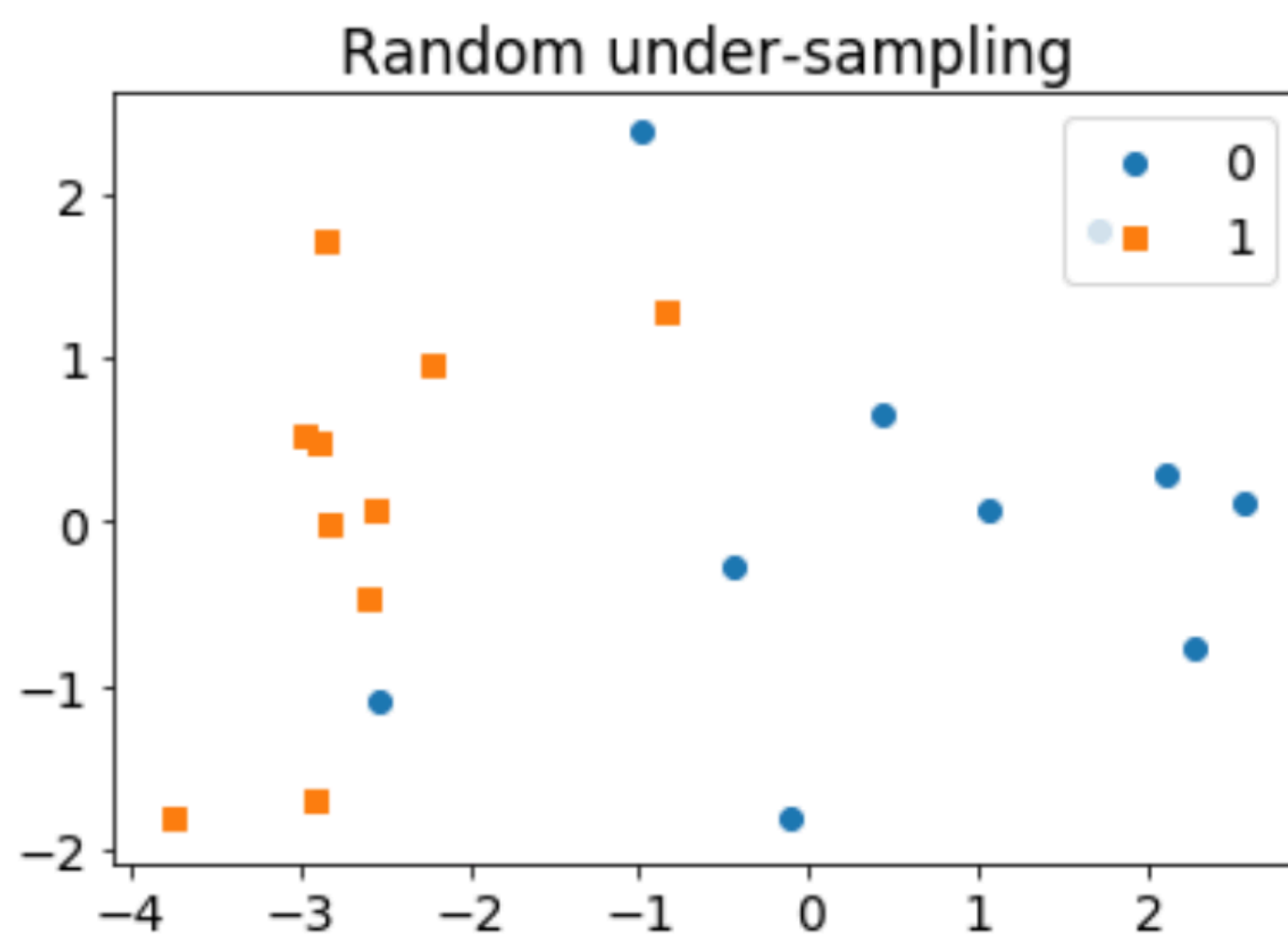
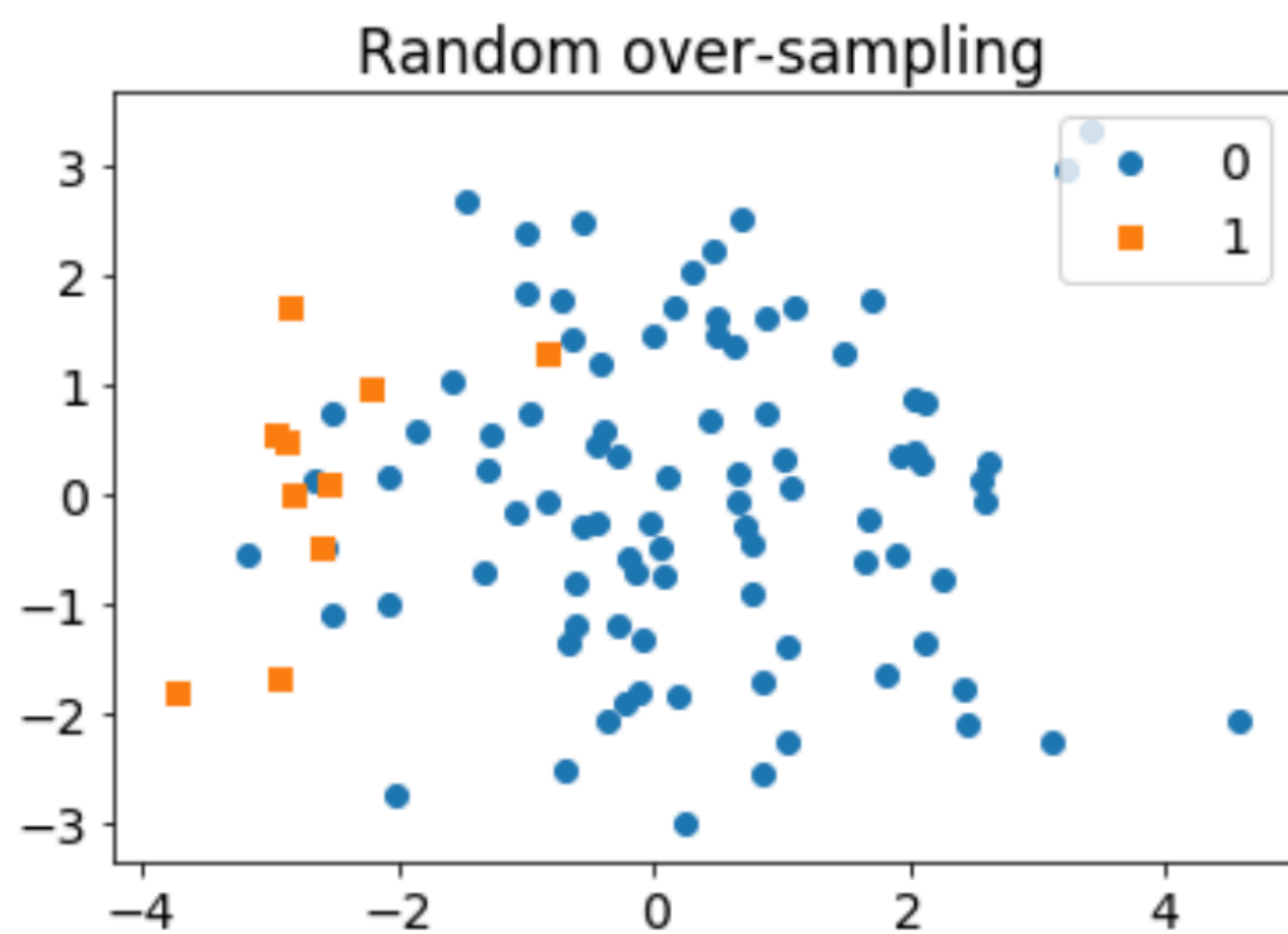
<http://www.statisticssolutions.com/sample-size-calculation-and-sample-size-justification-resampling/>

Hence to correct this bias in data, I used

- Resampling(Under-sampling, Over-sampling), Python imbalanced-learn module,
- Because the dataset has many dimensions (features) and our graphs will be 2D, we will reduce the size of the dataset using Principal Component Analysis (PCA)



- Random under-sampling and over-sampling with imbalanced-learn



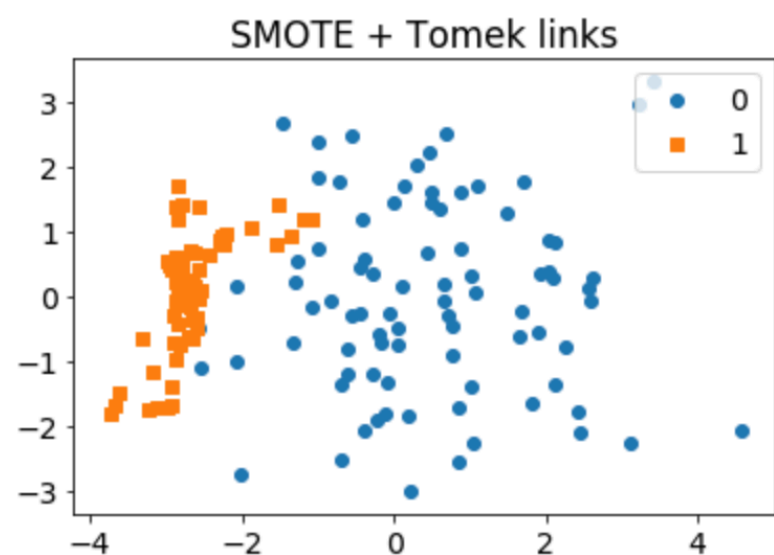
Over-sampling followed by under-sampling

Now, we will do a combination of over-sampling and under-sampling, using the SMOTE and Tomek links techniques:

```
: from imblearn.combine import SMOTETomek

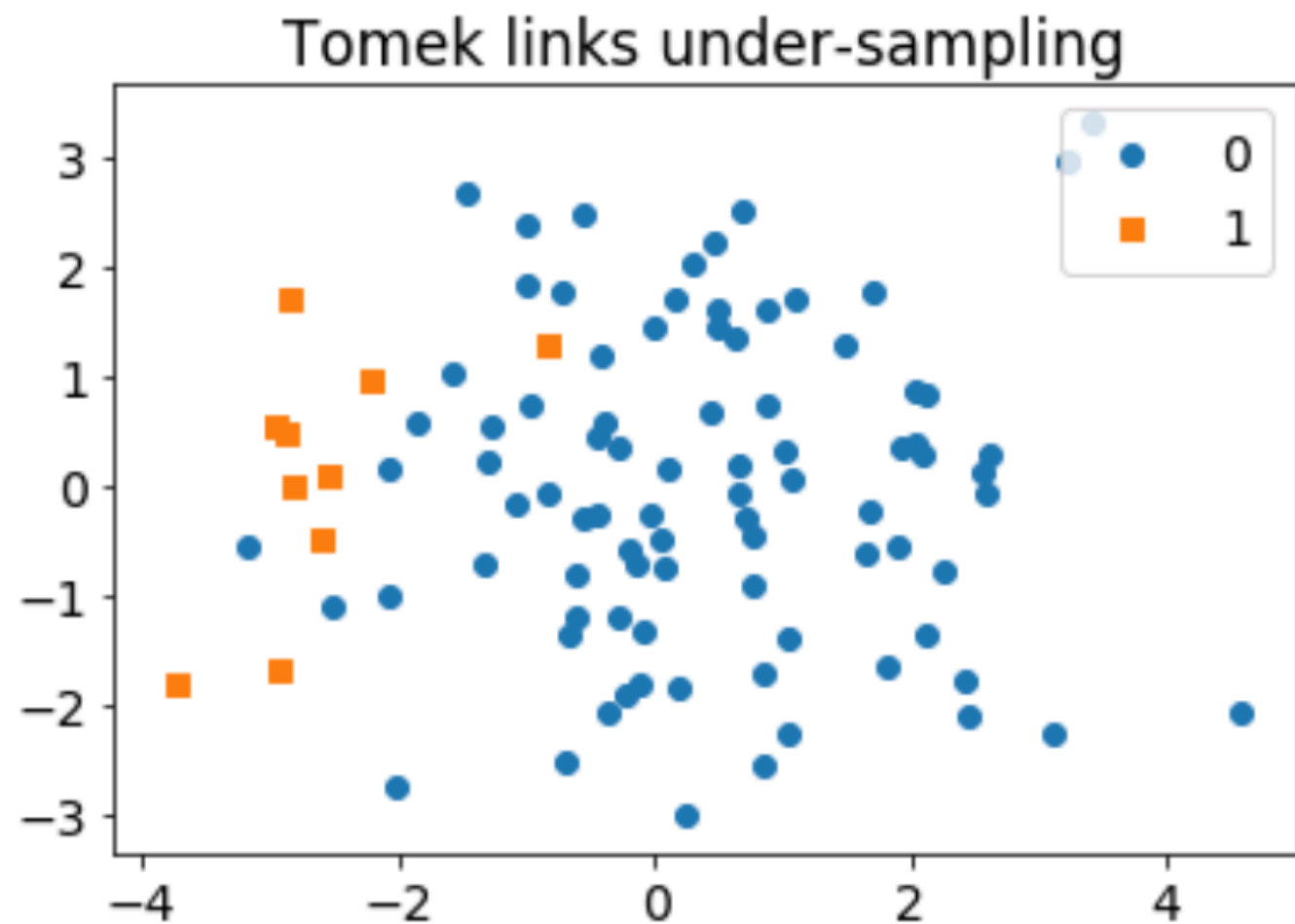
smt = SMOTETomek(ratio='auto')
X_smt, y_smt = smt.fit_sample(X, y)

plot_2d_space(X_smt, y_smt, 'SMOTE + Tomek links')
```

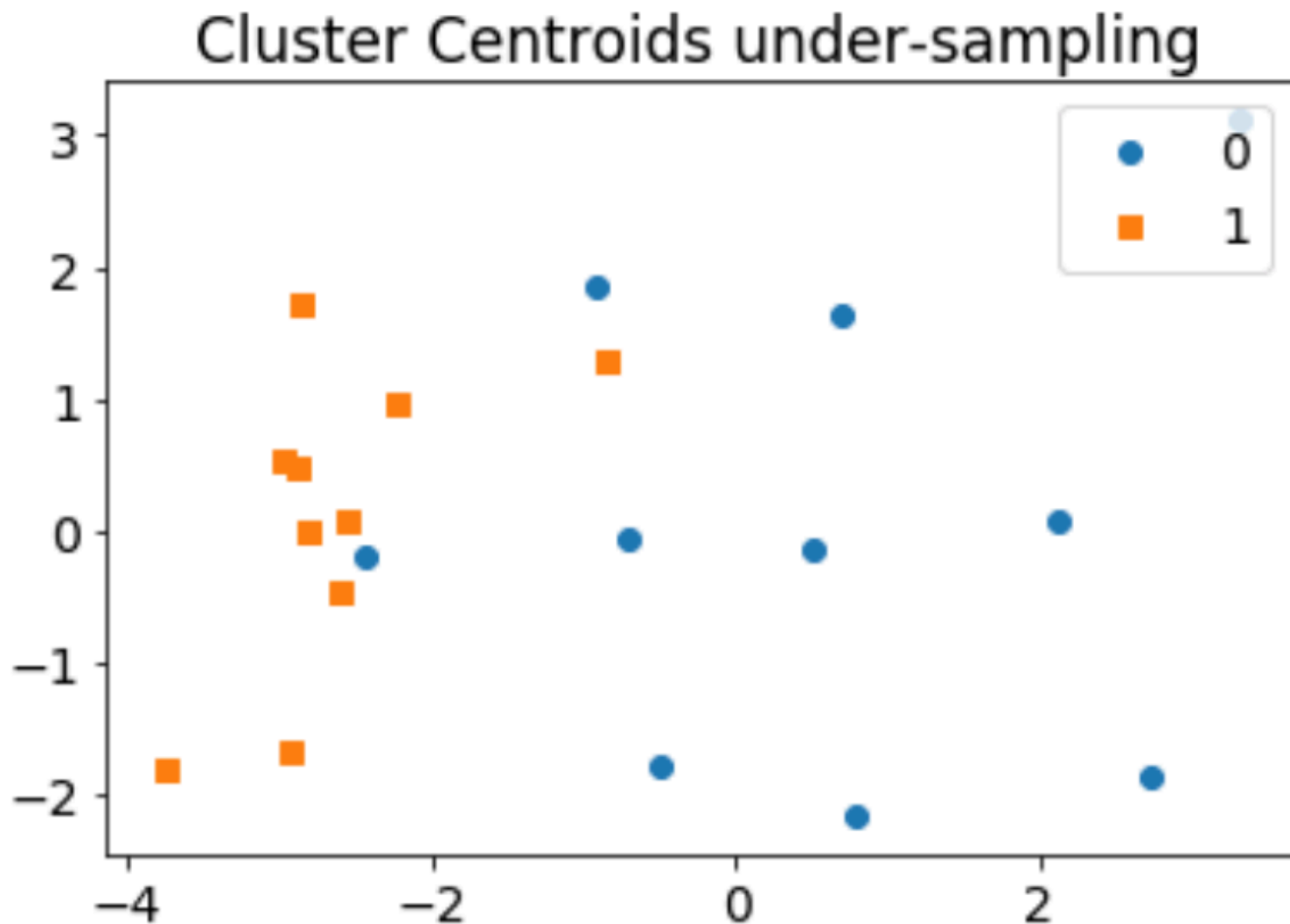


- Under-sampling: Tomek links remove unwanted overlap between classes where majority class links are removed until all minimally distanced nearest neighbor pairs are of the same class. courtesy:
https://en.wikipedia.org/wiki/Oversampling_and_undersampling_in_data_analysis

In the code below, we'll use ratio='majority' to resample the majority class.



- Under-sampling: Cluster centroids is a method that replaces cluster of samples by the cluster centroid of a K-means algorithm, where the number of clusters is set by the level of undersampling. courtesy:
https://en.wikipedia.org/wiki/Oversampling_and_undersampling_in_data_analysis



- Over-sampling: SMOTE

There are a number of methods available to oversample a dataset used in a typical classification problem (using a classification algorithm to classify a set of images, given a labelled training set of images). The most common technique is known as SMOTE: Synthetic Minority Over-sampling Technique.

We'll use `ratio='minority'` to resample the minority class.

- Over-sampling followed by under-sampling

Next I combine over-sampling and under-sampling, using the SMOTE and Tomek links techniques:

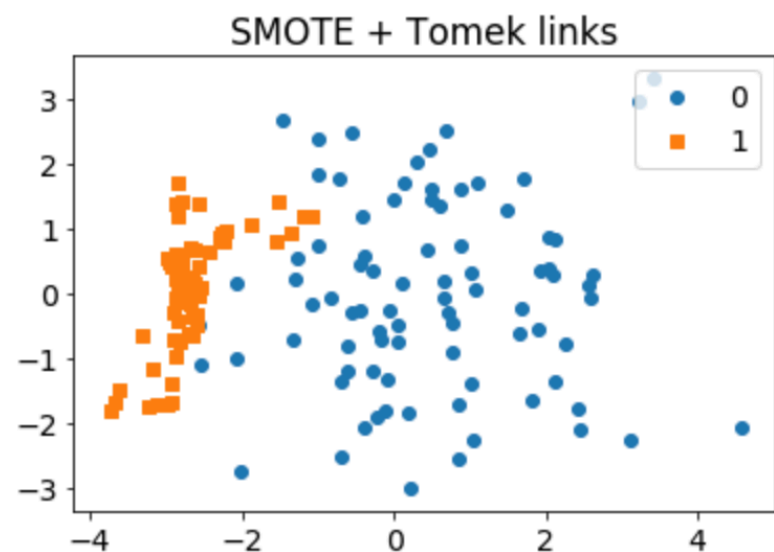
Over-sampling followed by under-sampling

Now, we will do a combination of over-sampling and under-sampling, using the SMOTE and Tomek links techniques:

```
: from imblearn.combine import SMOTETomek

smt = SMOTETomek(ratio='auto')
X_smt, y_smt = smt.fit_sample(X, y)

plot_2d_space(X_smt, y_smt, 'SMOTE + Tomek links')
```



Step 5: Evaluation

- Classification accuracy: percentage of correct predictions
- Precision attempts to answer the following question: What proportion of positive identifications was actually correct?

Precision is defined as follows: $Tp / Tp + Fp$

- Recall attempts to answer the following question: What proportion of actual positives was identified correctly?

Mathematically, recall is defined as follows: $Tp / Tp + Fn$

(Ref:

<https://developers.google.com/machine-learning/crash-course/classification/precision-and-recall>

- F1 score combines precision and recall relative to a specific positive class -The F1 score can be interpreted as a weighted average of the precision and recall, where an F1 score reaches its best value at 1 and worst at 0

Following are the Classification Reports Before & After Applying Imbalance in Datasets resolution for:

Logistic Regression

Before

	precision	recall	f1-score	support
0	0.92	1.00	0.96	56648
1	0.00	0.00	0.00	4855
avg / total	0.85	0.92	0.88	61503
0.9210607612636782				

After

	precision	recall	f1-score	support
0	1.00	0.79	0.88	14
1	0.87	1.00	0.93	20
avg / total	0.92	0.91	0.91	34
0.9117647058823529				

Based on above reports, Imbalance in dataset issue resolution does improve significantly:

- Precision
- F1 score

While accuracy & recall almost remain same. Hence Imbalance issue resolution does improve data model performance.

Overall Evaluation Report:

	Logistic Regression	Random Forest Classification	Decision Tree Classification
ACCURACY SCORE	0.82	0.91	0.96
PRECISION	0.84	0.91	0.96
RECALL	0.82	0.91	0.96
F1 SCORE	0.82	0.91	0.96
SUPPORT	34	43	51

Classification Reports for Additional Machine Learning Models:

Random Forest Classification

	precision	recall	f1-score	support
0	1.00	0.89	0.94	18
1	0.93	1.00	0.96	25
avg / total	0.96	0.95	0.95	43
0.9534883720930233				

Decision Tree Classification

	precision	recall	f1-score	support
0	1.00	1.00	1.00	23
1	1.00	1.00	1.00	28
avg / total	1.00	1.00	1.00	51

Step 6: Performance Tuning

Confusion Matrix is a table that describes the performance of a classification model:

- Every observation in the testing set is represented in exactly one box
- It's a 2x2 matrix because there are 2 response classes

(Courtesy:

<https://www.ritchieng.com/machine-learning-evaluate-classification-model/>)

Metrics computed from a confusion matrix are:

- Classification Accuracy: Overall, how often is the classifier correct
- Classification Error: Overall, how often is the classifier incorrect
- Sensitivity: When the actual value is positive, how often is the prediction correct
- Specificity: When the actual value is negative, how often is the prediction correct

Following are the Confusion Matrix outcomes for Logistic Regression, Random Forest and Decision Tree Classification Models respectively:

