

# Capstone Project II: Milestone Report

## Project Objective

## Project Question:

To predict the clients' repayment abilities. Also determine the factors or features that influence prediction by analyzing variety of alternative data--including telco and transactional information. The problem that the client wants to solve by throwing open data challenge in Kaggle is as follows:

- 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

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	GOODS_NAME	TYPE_NAME	NAME_INCO	NAME_EDUC	NAME_FAMI	NAME_HOU	REGION	POF	DAYS_BIRTH	DAYS_EM
100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5	351000	Unacompar	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 educ	Married	House / apar	0.003541	-16765	-11	
100004	0	Revolving loc	M	Y	Y	0	67500	135000	6750	135000	Unacompar	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	Unacompar	Working	Secondary / Civil marriag	House / apar	0.008019	-19005	-30		
100007	0	Cash loans	M	N	Y	0	121500	513000	21865.5	513000	Unacompar	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	Unacompar	Commercial	Higher educ	Married	House / apar	0.035792	-13778	-31	
100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075	1530000	Unacompar	State servant	Higher educ	Married	House / apar	0.003122	-18850	-4	
100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5	913500	Children	Pensioner	Secn/edu L	Married	House / apar	0.018634	-20099	3652	

Dataset given has 122 features with 307K +rows

Of the 122 features, following features showed **correlation** with the target variable 'TARGET'  
i.e. Loan Repayment Prediction:

- 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 after performing Data Wrangling steps as explained in google doc link given here:

[https://docs.google.com/document/d/1MesT5Q-V7EVpQ\\_Io5EbFUpeEIjFzgqaXq6oCLUjC3A/edit?usp=sharing](https://docs.google.com/document/d/1MesT5Q-V7EVpQ_Io5EbFUpeEIjFzgqaXq6oCLUjC3A/edit?usp=sharing)

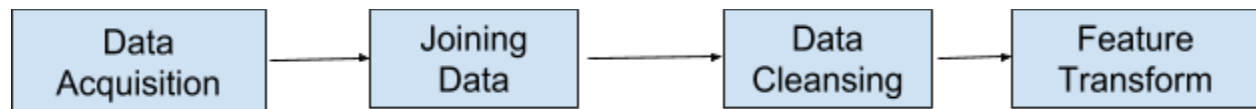
## Data Wrangling

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



The data is provided by Home Credit Group, a service dedicated to provided lines of credit (loans) to the unbanked population.

## Project Objective

Predicting whether or not a client will repay a loan or have difficulty is a critical business need, and Home Credit is hosting this competition on Kaggle to see what sort of models the machine learning community can develop to help them in this task.

application\_train/application\_test: 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.

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_CONTCODE	GEND	FLAG_OWN	FLAG_OWN_CNT	CHILD_AMT	INCOME_AMT	CREDIT_AMT	ANNUAL_AMT	GOOD_NAME	TYPE	NAME_INCOME	NAME_EDUC	NAME_FAMILY	NAME_HOU	REGION	POF	DAYS_BIRTH	DAYS_EM
100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5	351000	Unacompar	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 servan	Higher educ	Married	House / apar	0.003541	-16765	-11	
100004	0	Revolving loc	M	Y	Y	0	67500	135000	6750	135000	Unacompar	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	Unacompar	Working	Secondary /	Civil marriag	House / apar	0.008019	-19005	-30	
100007	0	Cash loans	M	N	Y	0	121500	513000	21865.5	513000	Unacompar	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 servan	Secondary /	Married	House / apar	0.035792	-16941	-15	
100009	0	Cash loans	F	Y	Y	1	171000	1560726	41301	1395000	Unacompar	Commercial	Higher educ	Married	House / apar	0.035792	-13778	-31	
100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075	1530000	Unacompar	State servan	Higher educ	Married	House / apar	0.003122	-18850	-4	
100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5	913500	Children	Pensioner	Secnary	Low L	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)

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
<b>COMMONAREA_MEDI</b>	214865	69.9
<b>COMMONAREA_AVG</b>	214865	69.9
<b>COMMONAREA_MODE</b>	214865	69.9
<b>NONLIVINGAPARTMENTS_MEDI</b>	213514	69.4
<b>NONLIVINGAPARTMENTS_MODE</b>	213514	69.4
<b>NONLIVINGAPARTMENTS_AVG</b>	213514	69.4
<b>FONDKAPREMONT_MODE</b>	210295	68.4
<b>LIVINGAPARTMENTS_MODE</b>	210199	68.4
<b>LIVINGAPARTMENTS_MEDI</b>	210199	68.4
<b>LIVINGAPARTMENTS_AVG</b>	210199	68.4
<b>FLOORSMIN_MODE</b>	208642	67.8
<b>FLOORSMIN_MEDI</b>	208642	67.8
<b>FLOORSMIN_AVG</b>	208642	67.8
<b>YEARS_BUILD_MODE</b>	204488	66.5
<b>YEARS_BUILD_MEDI</b>	204488	66.5
<b>YEARS_BUILD_AVG</b>	204488	66.5
<b>OWN_CAR_AGE</b>	202929	66.0
<b>LANDAREA_AVG</b>	182590	59.4
<b>LANDAREA_MEDI</b>	182590	59.4
<b>LANDAREA_MODE</b>	182590	59.4

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

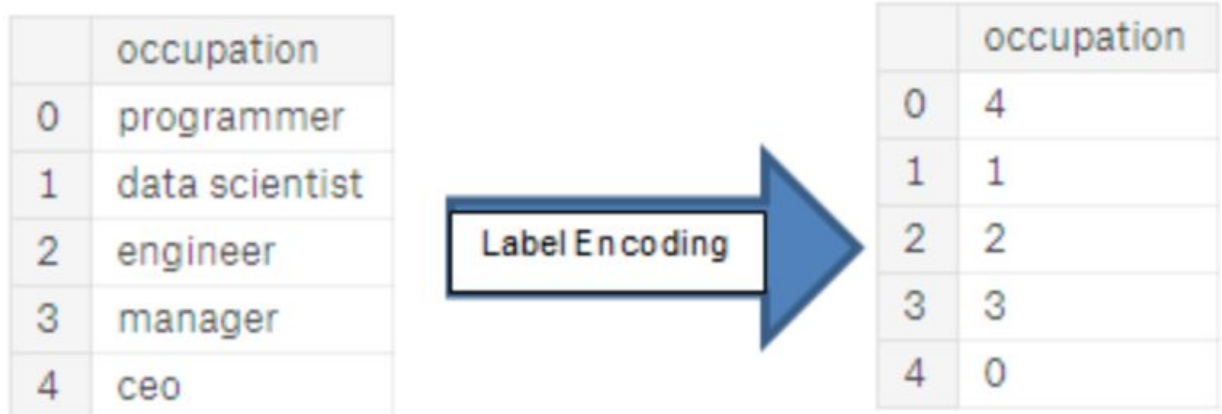
```
# Number of each type of column  
df_train.dtypes.value_counts()
```

```
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. A machine learning model unfortunately cannot deal with categorical variables (except for some models such as LightGBM). Therefore, we have to find a way to encode (represent) these variables as numbers before handing them off to the model. There are two main ways to carry out this process:

Label encoding: assign each unique category in a categorical variable with an integer. No new columns are created. An example is shown below image





One-hot encoding: create a new column for each unique category in a categorical variable. Each observation receives a 1 in the column for its corresponding category and a 0 in all other new columns.



	occupation
0	programmer
1	data scientist
2	engineer
3	manager
4	ceo

	occupation_ceo	occupation_data scientist	occupation_engineer	occupation_manager	occupation_programmer
0	0	0	0	0	1
1	0	1	0	0	0
2	0	0	1	0	0
3	0	0	0	1	0
4	1	0	0	0	0

The problem with label encoding is that it gives the categories an arbitrary ordering. The value assigned to each of the categories is random and does not reflect any inherent aspect of the category. In the example above, programmer receives a 4 and data scientist a 1, but if we did the same process again, the labels could be reversed or completely different. The actual assignment of the integers is arbitrary. Therefore, when we perform label encoding, the model might use the relative value of the feature (for example programmer = 4 and data scientist = 1) to assign weights which is not what we want. If we only have two unique values for a categorical variable (such as Male/Female), then label encoding is fine, but for more than 2 unique categories, one-hot encoding is the safe option.

There is some debate about the relative merits of these approaches, and some models can deal with label encoded categorical variables with no issues. Here is a good [Stack Overflow discussion](#). I think (and this is just a personal opinion) for categorical variables with many classes, one-hot encoding is the safest approach because it does not impose arbitrary values to categories. The only downside to one-hot encoding is that the number of features (dimensions of the data) can explode with categorical variables with many categories. To deal with this, we can perform one-hot encoding followed by PCA or other dimensionality reduction methods to reduce the number of dimensions (while still trying to preserve information).

In this notebook, we will use Label Encoding for any categorical variables with only 2 categories and One-Hot Encoding for any categorical variables with more than 2 categories. This process may need to change as we get further into the project, but for now, we will see where this gets us. (We will also not use any dimensionality reduction in this notebook but will explore in future iter

Let's implement the policy described above: for any categorical variable (`dtype == object`) with 2 unique categories, we will use label encoding, and for any categorical variable with more than 2 unique categories, we will use one-hot encoding.

For label encoding, we use the Scikit-Learn `LabelEncoder` and for one-hot encoding, the pandas `get_dummies(df)` function.

3 columns were label encoded.

## VI. Aligning Training and Testing Data

There need to be the same features (columns) in both the training and testing data. One-hot encoding has created more columns in the training data because there were some categorical variables with categories not represented in the testing data. To remove the columns in the training data that are not in the testing data, we need to align the dataframes. First we extract the target column from the training data (because this is not in the testing data but we need to keep this information). When we do the align, we must make sure to set `axis = 1` to align the dataframes based on the columns and not on the rows!

('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. At some point we probably will want to try dimensionality reduction (removing features that are not relevant) to reduce the size of the datasets.

## VII. Missing Data Strategy / Outlier Analysis

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

One problem we always want to be on the lookout for when doing EDA is anomalies within the data. These may be due to mis-typed numbers, errors in measuring equipment, or they could be valid but extreme measurements. One way to support anomalies quantitatively is by looking at the statistics of a column using the describe method. The numbers in the DAYS\_BIRTH column are negative because they are recorded relative to the current loan application. To see these stats in years, we can mutiple by -1 and divide by the number of days in a year:

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

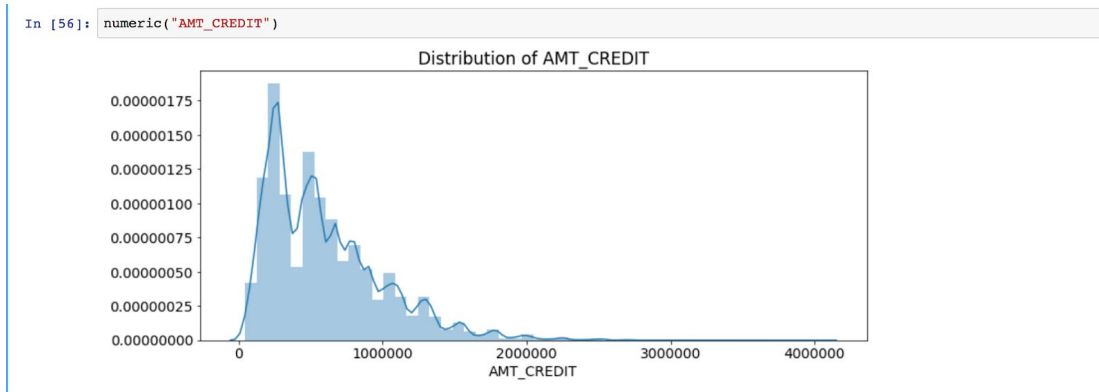
```
count      307511.000000
mean         43.936973
std          11.956133
min          20.517808
25%          34.008219
50%          43.150685
75%          53.923288
max          69.120548
Name: DAYS_BIRTH, dtype: float64
```

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

## Exploratory Data Analysis

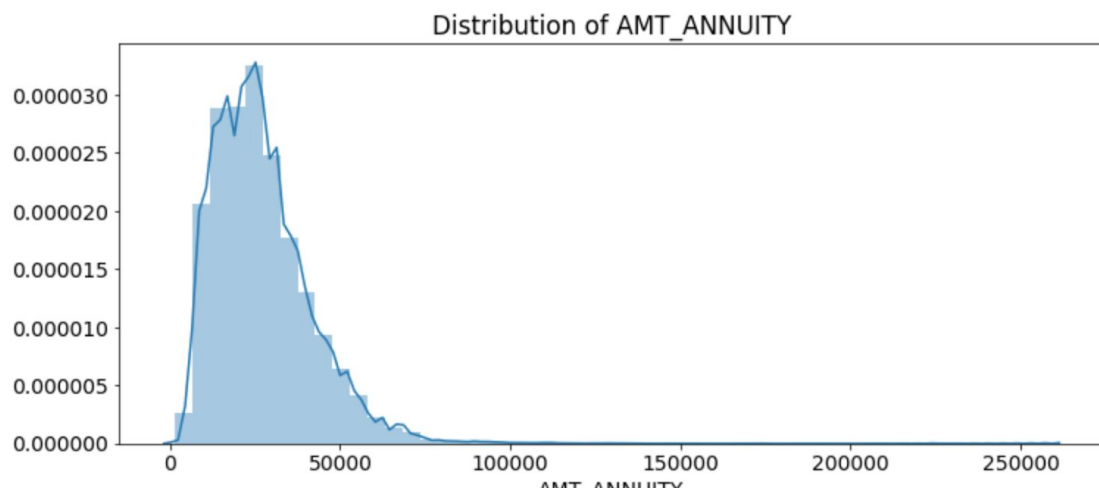
Visualized distribution of numeric independent features

- AMT\_CREDIT



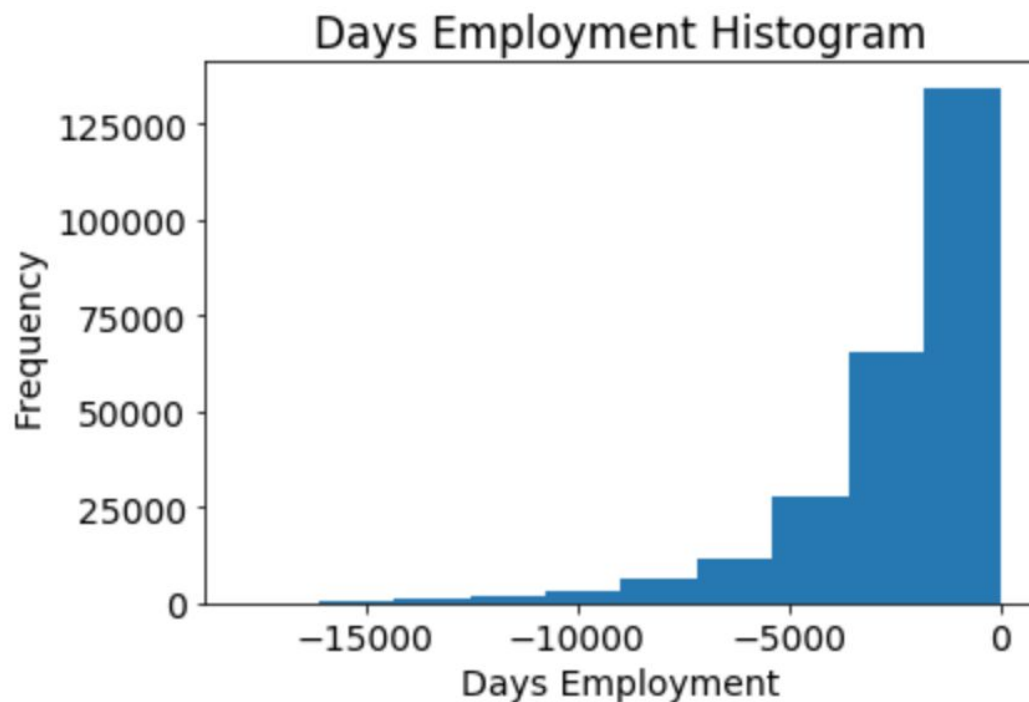
- AMT\_ANNUIITY

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

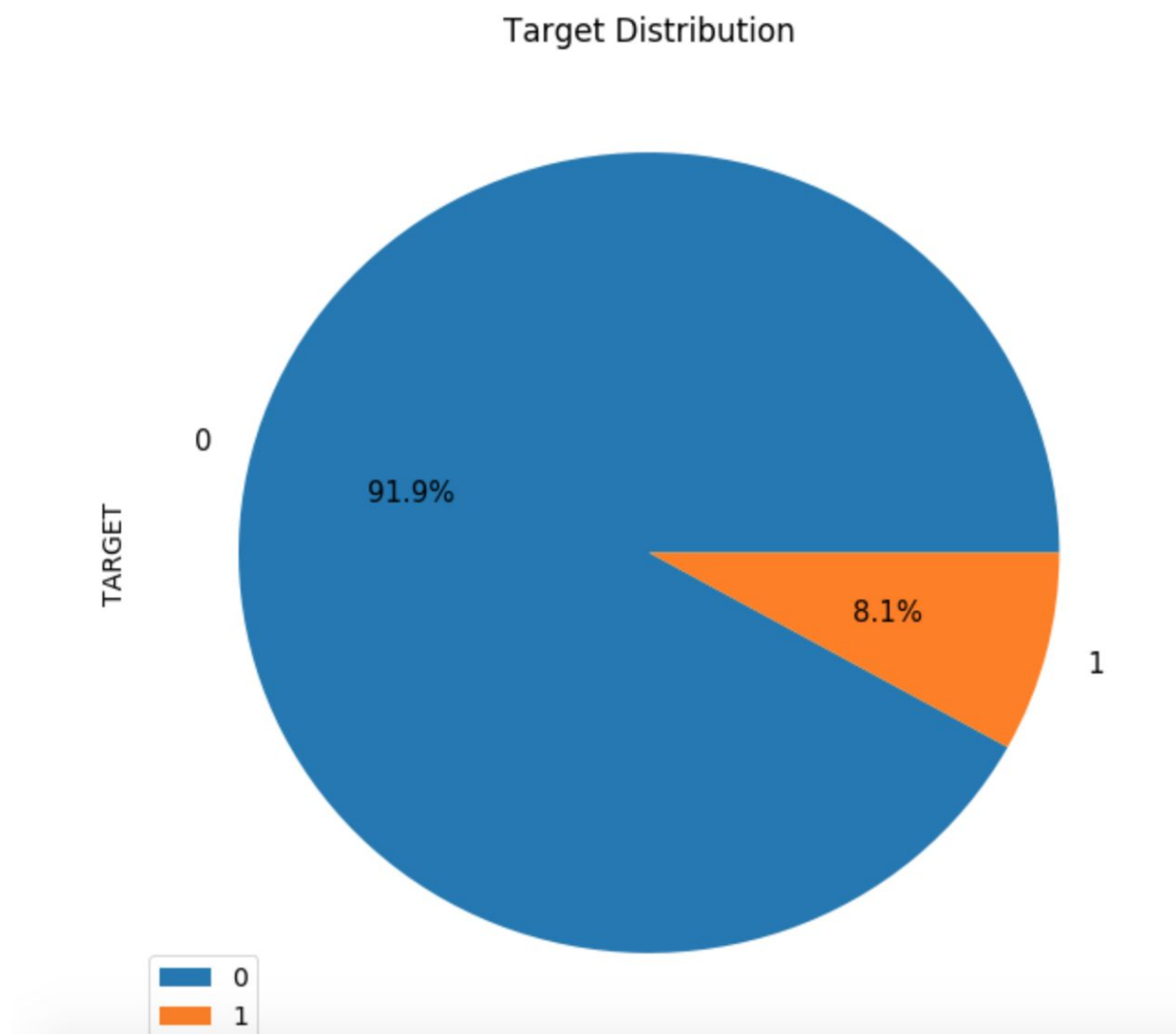


2. The 'Days Employment' distribution looks to be much more in line with what we would expect, and we also have created a new column to tell the model that these values were originally anomalous (because we will have to fill in the nans with some value, probably the median of the column). The other columns 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.

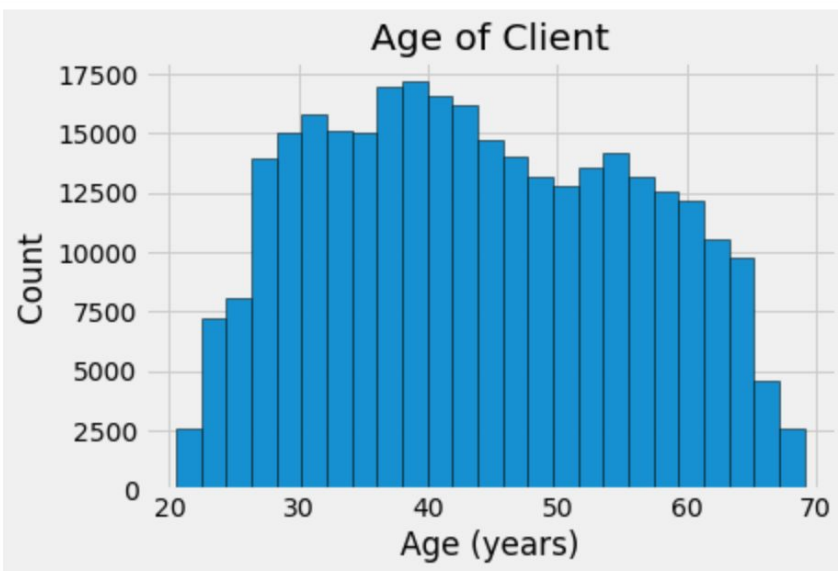
#### 4. Effect of Age on Repayment

By itself, the distribution of age does not tell us much other than that there are no outliers as all the ages are reasonable. 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.

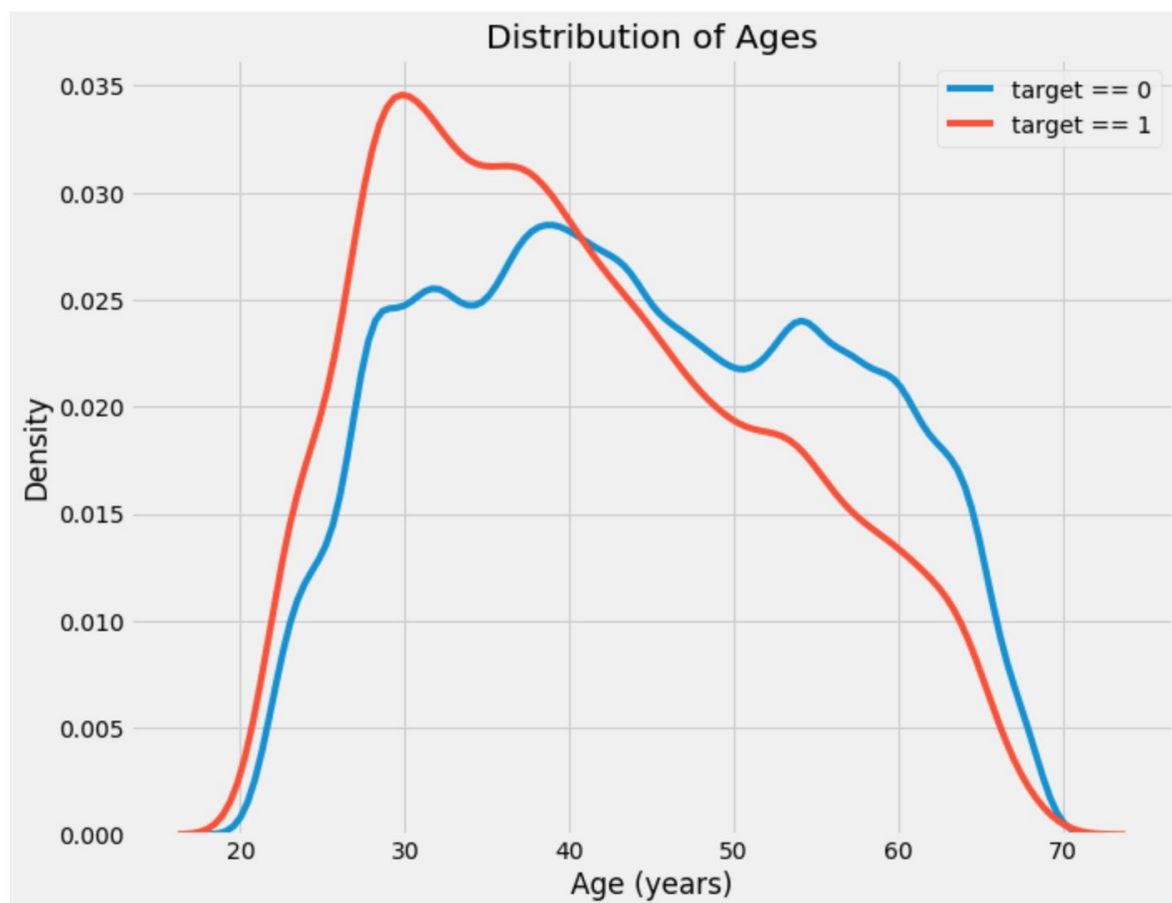
```
# Set the style of plots
plt.style.use('fivethirtyeight')

# 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

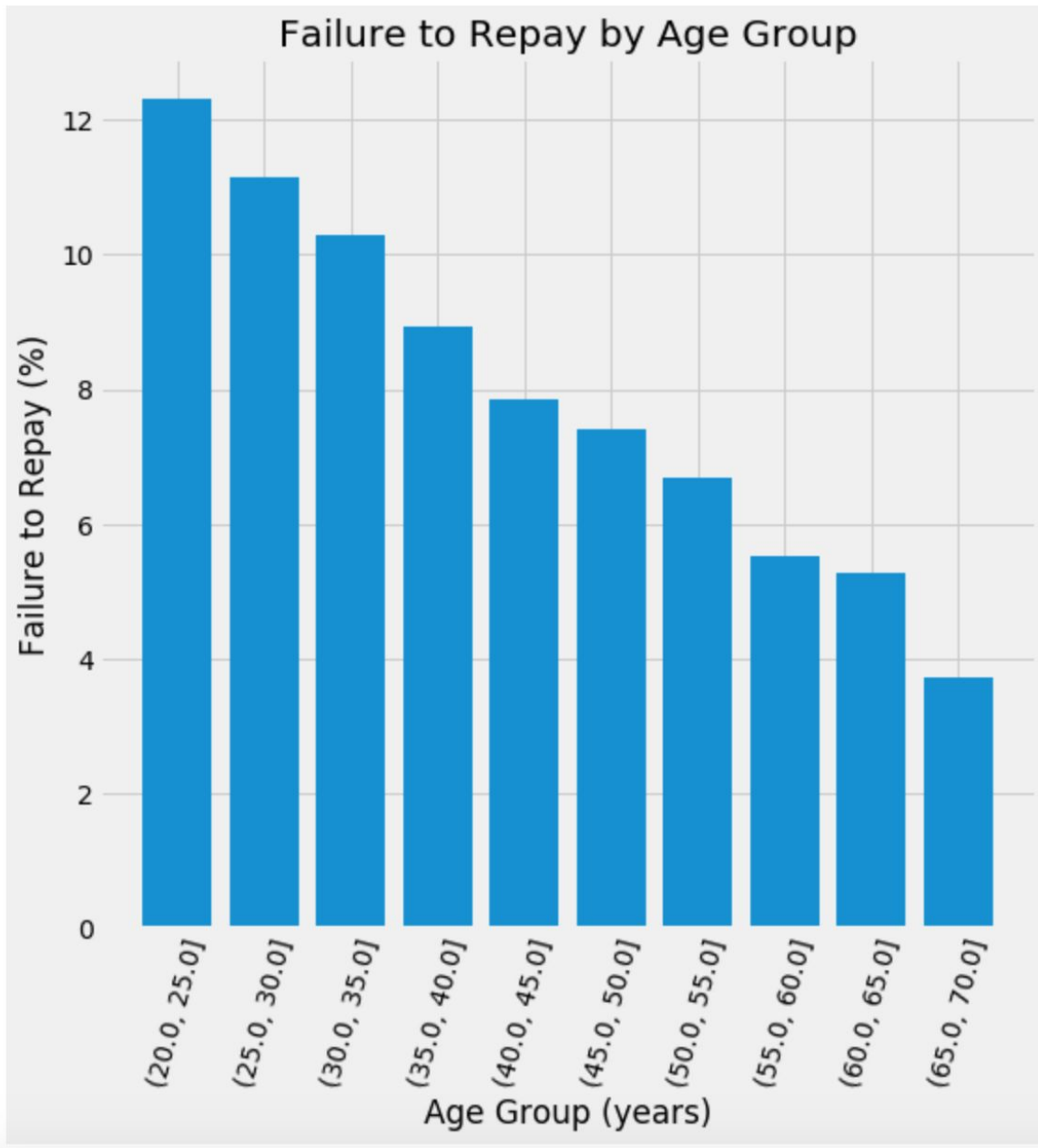
	TARGET	DAYS_BIRTH	YEARS_BIRTH
YEARS_BINNED			
(20.0, 25.0]	0.123036	8532.795625	23.377522
(25.0, 30.0]	0.111436	10155.219250	27.822518
(30.0, 35.0]	0.102814	11854.848377	32.479037
(35.0, 40.0]	0.089414	13707.908253	37.555913
(40.0, 45.0]	0.078491	15497.661233	42.459346
(45.0, 50.0]	0.074171	17323.900441	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.



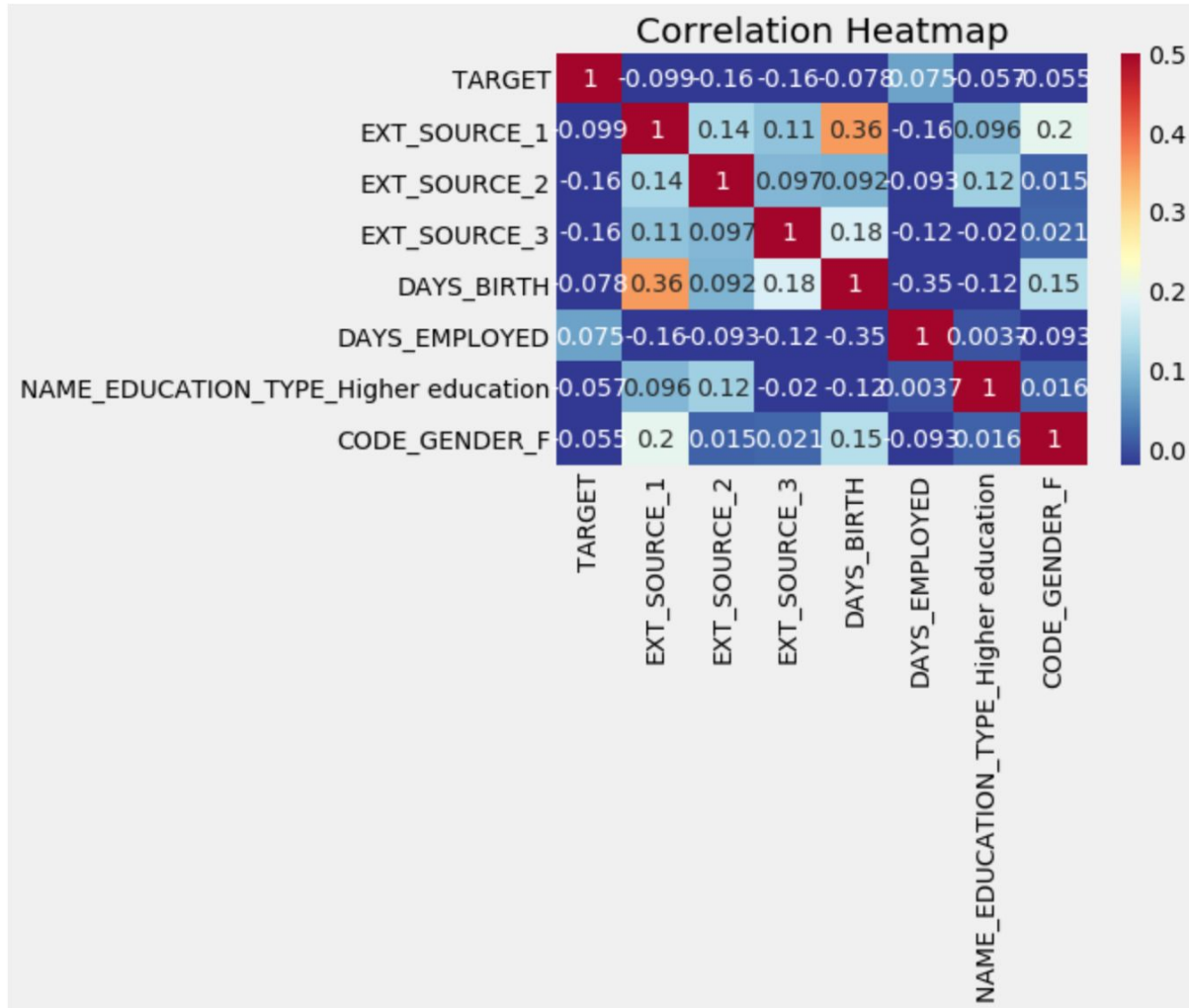


## 8. Exterior Sources

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.

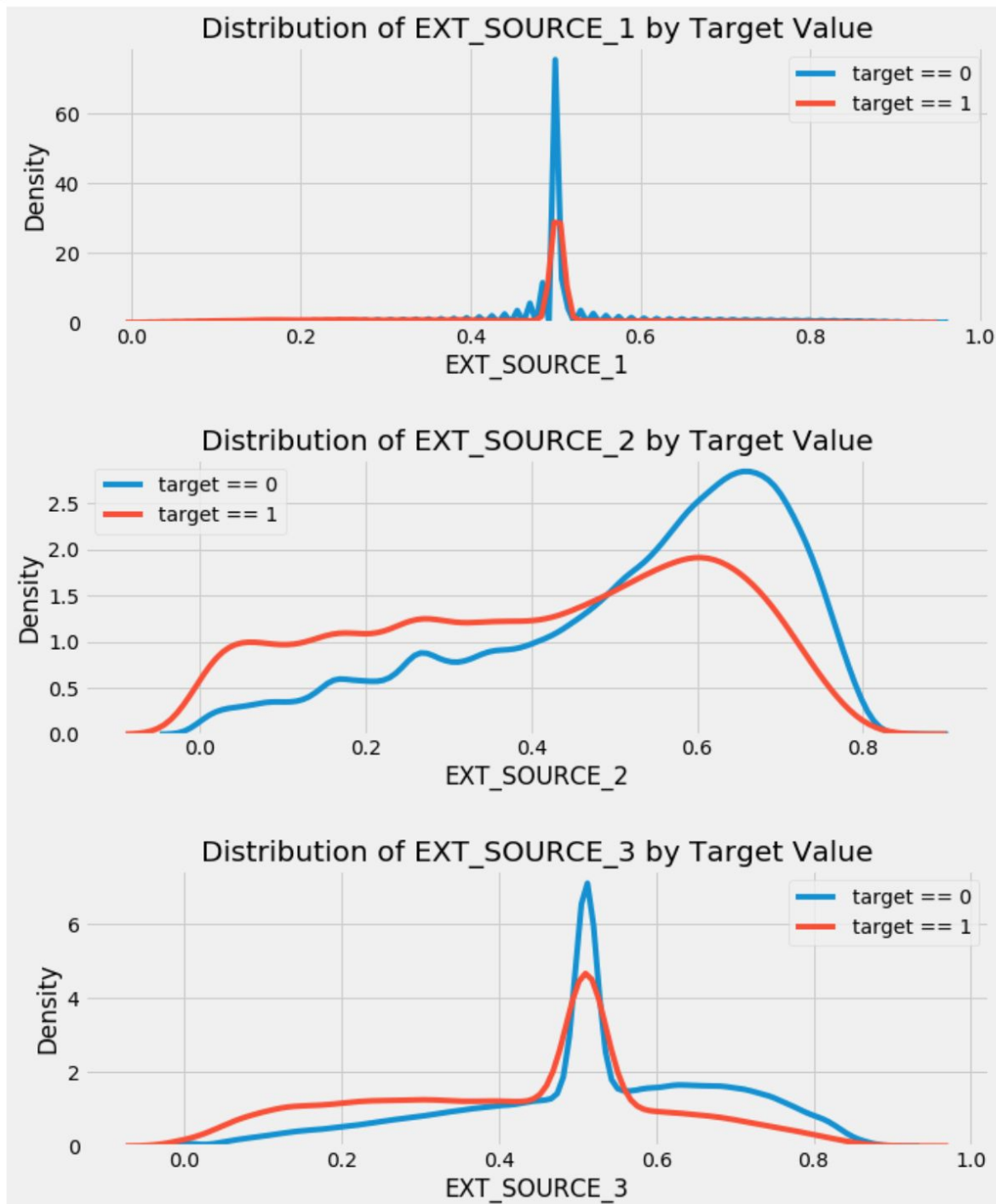
First, we can show the correlations of the EXT\_SOURCE features with the target and with each other.



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

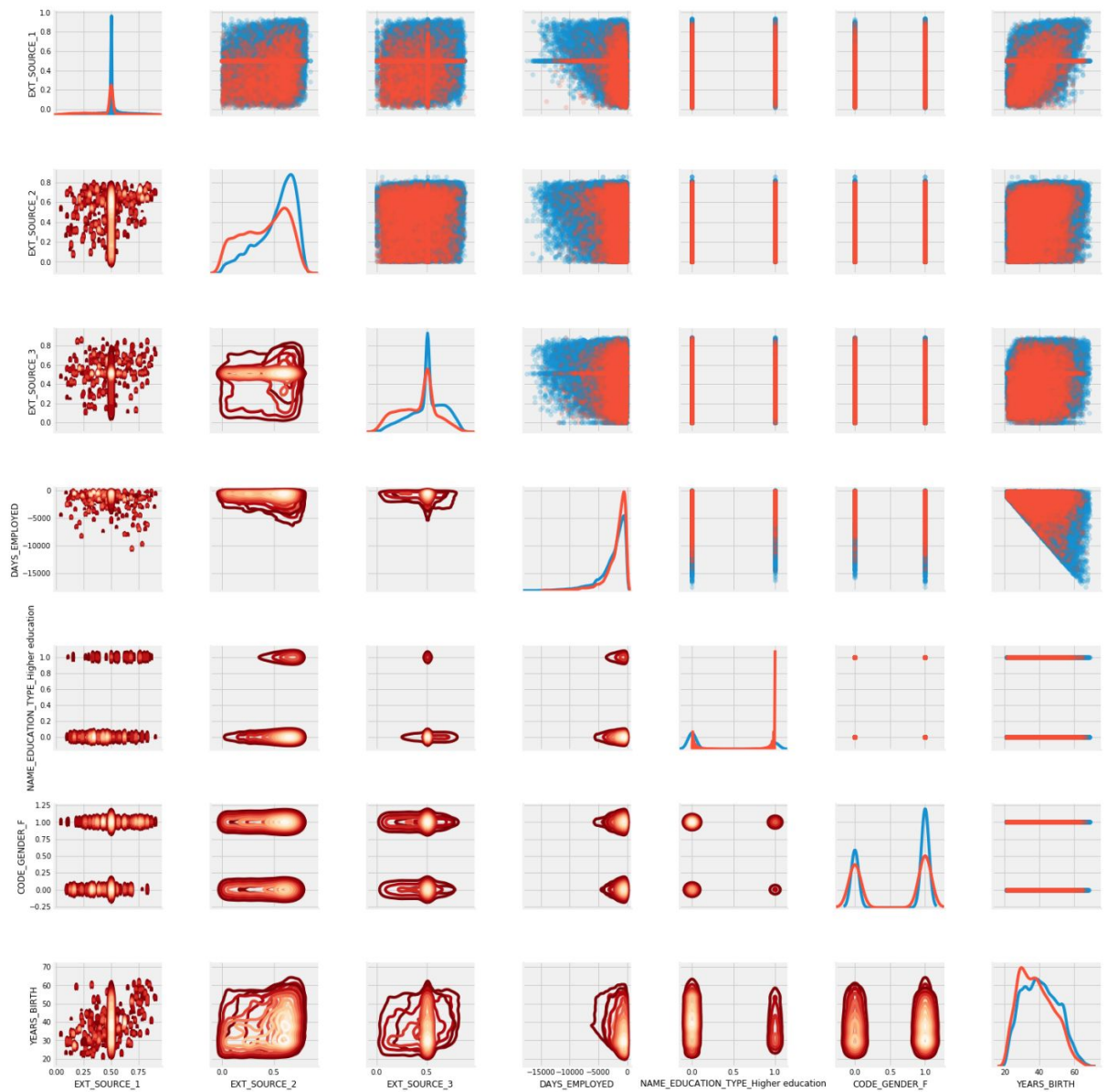
- Employment Duration



## 1. Pairs Plot

As a final exploratory plot, we can make a pairs plot of the EXT\_SOURCE variables and the DAYS\_BIRTH variable. The Pairs Plot is a great exploration tool because it lets us see relationships between multiple pairs of variables as well as distributions of single variables. Here we are using the seaborn visualization library and the PairGrid function to create a Pairs Plot with scatterplots on the upper triangle, histograms on the diagonal, and 2D kernel density plots and correlation coefficients on the lower triangle.

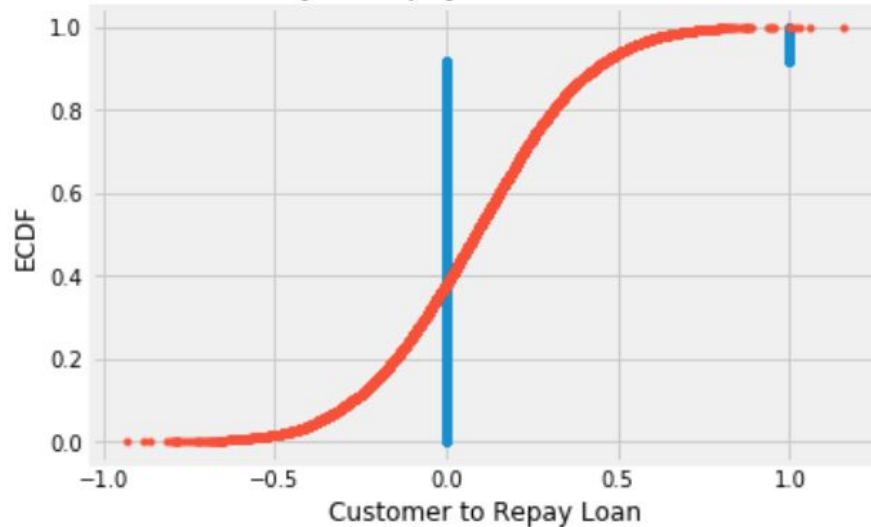
Ext Source and Age Features Pairs Plot



In this plot, the red indicates loans that were not repaid and the blue are loans that are paid. We can see the different relationships within the data. There does appear to be a moderate positive linear relationship between the EXT\_SOURCE\_1 and the DAYS\_BIRTH (or equivalently YEARS\_BIRTH), indicating that this feature may take into account the age of the client.

## 2. Empirical Continuous Distribution ECDF plot for Credit Loan Default Risk

Customer Credibility to Repay Loan VS Theoretical Normal Dist



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)

## 3. Checked Variance, Covariance, Standard Deviation, Pearson Correlation Coefficient

## Variance

```
: np.var(df_train[ 'TARGET' ])  
:  
: 0.0742116771655796
```

## Standard Deviation

```
: np.std(df_train[ 'TARGET' ])  
:  
: 0.27241820270602257
```

## Covariance

```
: np.cov(df_train[ 'TARGET' ], df_train[ 'EXT_SOURCE_1' ])  
:  
: array([[ 0.07421192, -0.0037652 ],  
        [-0.0037652 ,  0.01943096]])
```

```
: np.cov(df_train[ 'TARGET' ], df_train[ 'DAYS_BIRTH' ])  
:  
: array([[ 7.42119185e-02, -9.30133834e+01],  
        [-9.30133834e+01,  1.90443968e+07]])
```

## Pearson Correlation Coefficient

```
: np.corrcoef(df_train[ 'TARGET' ], df_train[ 'DAYS_BIRTH' ])  
:  
: array([[ 1.          , -0.07823931],  
        [-0.07823931,  1.          ]])
```

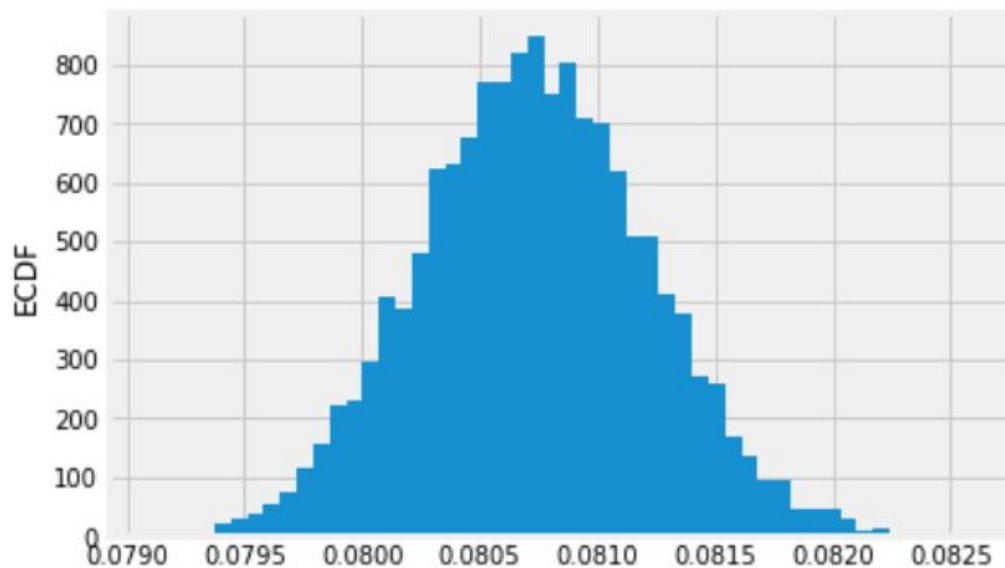
# 4. Confidence Interval

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

```
0.0004912536560492155  
0.0004897216997795728
```



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.

## 5. Hypothesis Testing

Hypothesis Testing

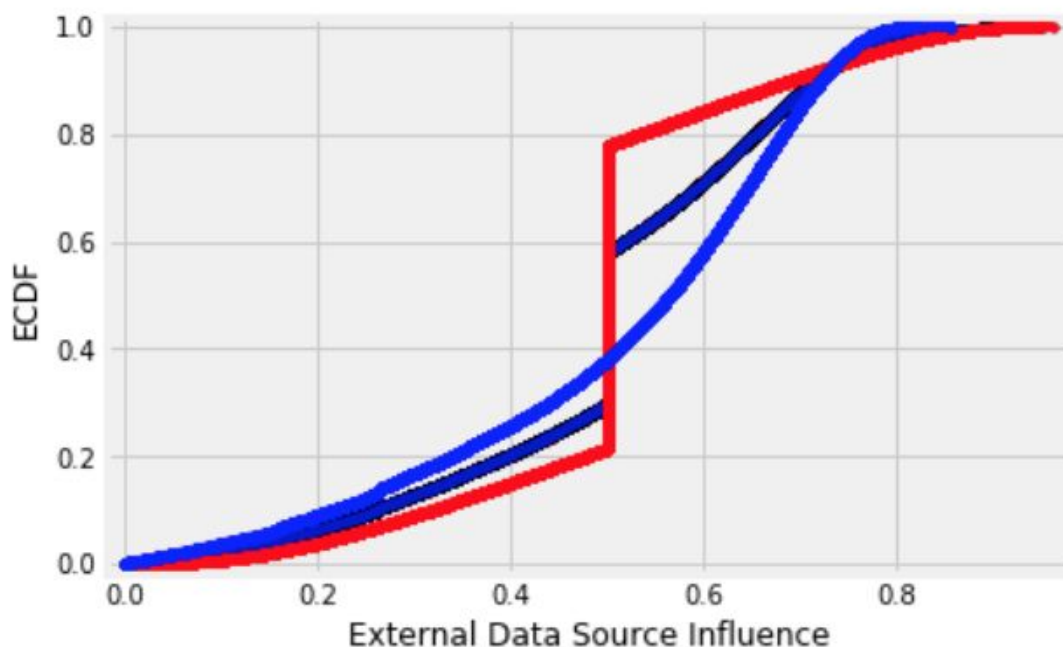


**Null Hypothesis- There is no significant difference between EXT\_SOURCE\_1 and EXT\_SOURCE\_2 mean on 'Ability to Repay Loan'**

**H0:  $\mu_{\text{EXT\_SOURCE\_1}} - \mu_{\text{EXT\_SOURCE\_2}} = 0$  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'**

**HA :  $\mu_{\text{EXT\_SOURCE\_1}} - \mu_{\text{EXT\_SOURCE\_2}} \neq 0$**



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.

Here is the link to iPython Notebook:

<https://github.com/rashi-n/Machine-Learning-Projects/blob/master/Capstone%20Projects/CS%201%20-%20EDA%20%26%20Inferential%20Statistics.ipynb>

6. Based on above EDA & Inferential analysis, I next propose to do in-depth analysis on the given dataset using one of the Supervised Machine Learning algorithms of type Classification since output datasets are provided and I can use this to predict future outcomes of target variable. Additionally this is Classification problem as dependent variable 'TARGET' is Boolean value indicating if given transaction of a Customer is at 'Credit Loan Default Risk' (value of 1) or not(value of 0).

Here is the link to iPython Notebook:

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