

# HOME CREDIT DEFAULT RISK PROJECT

## Project Aim:

Building a model to find out how capable each loan applicant is of repaying a loan, so that approving loans only for the applicants who are likely to repay the loan.

## Data Description and Overview:

There are 7 different sources of data:

1. application\_train/application\_test: The main training 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.
2. bureau: In this dataset it consists of data concerning client's previous credits from other financial institutions. Each previous credit has its own row in bureau, but one loan in the application data can have multiple previous credits.
3. bureau\_balance: It consists of monthly data about the previous credits in bureau. Each row is one month of a previous credit, and a single previous credit can have multiple rows, one for each month of the credit length.
4. previous\_application: The data of previous applications for loans at Home Credit of clients who have loans in the application data. Each current loan in the application data can have multiple previous loans. Each previous application has one row and is identified.
5. POS\_CASH\_BALANCE: It consists of monthly data about previous point of sale or cash loans clients have had with Home Credit. Each row is one month of a previous point of sale or cash loan, and a single previous loan can have many rows.
6. credit\_card\_balance: The monthly data about previous credit cards clients have had with Home Credit. Each row is one month of a credit card balance, and a single credit card can have many rows.
7. installments\_payment: The data of payment history for previous loans at Home Credit. There is one row for every made payment and one row for every missed payment. <https://medium.com/@praveenkotha/home-credit-default-risk-end-to-end-machine-learning-project-1871f52e3ef2>  
(<https://medium.com/@praveenkotha/home-credit-default-risk-end-to-end-machine-learning-project-1871f52e3ef2>)

## Project Outline

1. Exploratory Data Analysis: This section includes performing initial investigations on data so as to discover patterns with the help of summary statistics and graphical representations.
2. Feature Generation and Elimination: This section includes generating features using the main data source and combining other supplementary sources provided, after investigation and understanding the data. This step is applied on both train and test sets.
3. Model Application: This section includes applying a machine learning algorithm on the final train set to predict test set labels.

## 1.Exploratory Data Analysis:

## Importing Libraries

In [1]:

```
1 import pandas as pd
2 import sklearn
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import os
6 import warnings
7 import seaborn as sns
8 from sklearn.preprocessing import OneHotEncoder
9 from sklearn.impute import SimpleImputer
10 from sklearn.pipeline import Pipeline
11 from sklearn.compose import ColumnTransformer
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.metrics import roc_auc_score
14 from sklearn.calibration import CalibratedClassifierCV
15 from sklearn.metrics import confusion_matrix
16 from sklearn.ensemble import RandomForestClassifier
17 from sklearn.metrics import accuracy_score
18 import plotly.offline as py
19 import plotly.graph_objs as go
20 from plotly.offline import init_notebook_mode, iplot
21 from sklearn.model_selection import train_test_split
22 init_notebook_mode(connected=True)
23 import pickle
24 import gc
25 from plotly import tools
26
27 warnings.filterwarnings('ignore')
28 %matplotlib inline
```

Importing Data: Reading all 7 files into python

In [2]:

```
1 directory = 'C:\\CreditRiskProject\\Datasets'
2 train_set = pd.read_csv(directory+ '\\application_train.csv')
3 test_set = pd.read_csv(directory+ '\\application_test.csv')
4 bureau = pd.read_csv(directory+ '\\bureau.csv')
5 pos_cash = pd.read_csv(directory+ '\\POS_CASH_balance.csv')
6 previous_applicaton = pd.read_csv(directory+ '\\previous_application.csv')
7 credit_card_balance = pd.read_csv(directory+ '\\credit_card_balance.csv')
8 installments_payments = pd.read_csv(directory+ '\\installments_payments.csv')
```

We can see how the main dataset looks like

In [3]:

```
1 train_set.head()
```

Out[3]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY
0	100002	1	Cash loans	M	N	N
1	100003	0	Cash loans	F	N	N
2	100004	0	Revolving loans	M	Y	N
3	100006	0	Cash loans	F	N	N
4	100007	0	Cash loans	M	N	N

5 rows × 122 columns

We can check the statistical values of numeric columns. Although 'SK\_ID\_CURR' and 'TARGET' are numeric columns, we should not include them as SK\_ID\_CURR is the unique id and TARGET is the binary classification label.

In [19]:

```
1 numerics = train_set._get_numeric_data()
2 numerics = numerics.drop(columns=['SK_ID_CURR', 'TARGET'])
3 numerics.describe()
```

Out[19]:

	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE
count	307511.000000	3.075110e+05	3.075110e+05	307499.000000	3.072330e+0
mean	0.417052	1.687979e+05	5.990260e+05	27108.573909	5.383962e+0
std	0.722121	2.371231e+05	4.024908e+05	14493.737315	3.694465e+0
min	0.000000	2.565000e+04	4.500000e+04	1615.500000	4.050000e+0
25%	0.000000	1.125000e+05	2.700000e+05	16524.000000	2.385000e+0
50%	0.000000	1.471500e+05	5.135310e+05	24903.000000	4.500000e+0
75%	1.000000	2.025000e+05	8.086500e+05	34596.000000	6.795000e+0
max	19.000000	1.170000e+08	4.050000e+06	258025.500000	4.050000e+0

8 rows × 104 columns

We can see the ratio of null values for each column. Columns having more null values than the threshold should be checked before starting the model. Since these features might be misleading to the model, we can consider removing them later on.

In [20]:

```

1 Desc = train_set.describe().transpose().reset_index()
2 Desc = Desc.rename(columns = {"index":"Features"})
3 Desc["NullCount"] = 0
4 Desc["NullRatio"] = 0
5 for i in train_set.columns:
6     Null = train_set[i].isna().sum()
7     Ratio = float((train_set[i].isna().sum())) / float(len(train_set))
8     Desc.loc[Desc['Features'] == i, ["NullCount"]] = Null
9     Desc.loc[Desc['Features'] == i, ["NullRatio"]] = Ratio
10
11 NullsTable = Desc[['Features', 'count', 'NullCount', 'NullRatio']].sort_values(by = ['NullRatio'])
12 NullsTable.head(10)

```

Out[20]:

	Features	count	NullCount	NullRatio
50	COMMONAREA_MODE	92646.0	214865	0.698723
36	COMMONAREA_AVG	92646.0	214865	0.698723
64	COMMONAREA_MEDI	92646.0	214865	0.698723
44	NONLIVINGAPARTMENTS_AVG	93997.0	213514	0.694330
58	NONLIVINGAPARTMENTS_MODE	93997.0	213514	0.694330
72	NONLIVINGAPARTMENTS_MEDI	93997.0	213514	0.694330
56	LIVINGAPARTMENTS_MODE	97312.0	210199	0.683550
42	LIVINGAPARTMENTS_AVG	97312.0	210199	0.683550
70	LIVINGAPARTMENTS_MEDI	97312.0	210199	0.683550
68	FLOORSMIN_MEDI	98869.0	208642	0.678486

As we can see, many columns have high ratio of missing values. We should not use these columns having null values over 10% for our model, and deal with the null values of other columns with less than 10%.

## Distribution Of Data Among Positive and Negative Classes

We should check whether the train dataset is balanced or imbalanced.

In [4]:

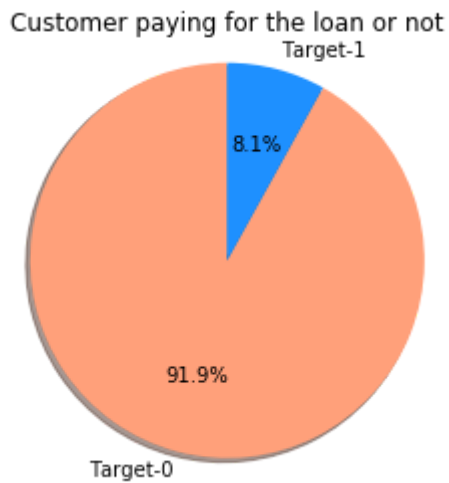
```

1 train_set_Target1 = train_set[train_set['TARGET'] == 1]
2 train_set_Target0 = train_set[train_set['TARGET'] == 0]

```

In [5]:

```
1 temp = train_set["TARGET"].value_counts()
2 temp2 = pd.DataFrame({'labels': ["Target-0", "Target-1"],
3                             'values': temp.values
4                             })
5 colors = ['lightsalmon', 'dodgerblue']
6 plt.pie(temp2["values"], labels=temp2["labels"], colors=colors,
7         autopct='%1.1f%%', shadow=True, startangle=90)
8 plt.title("Customer paying for the loan or not")
9
10 plt.axis('equal')
11 plt.show()
```



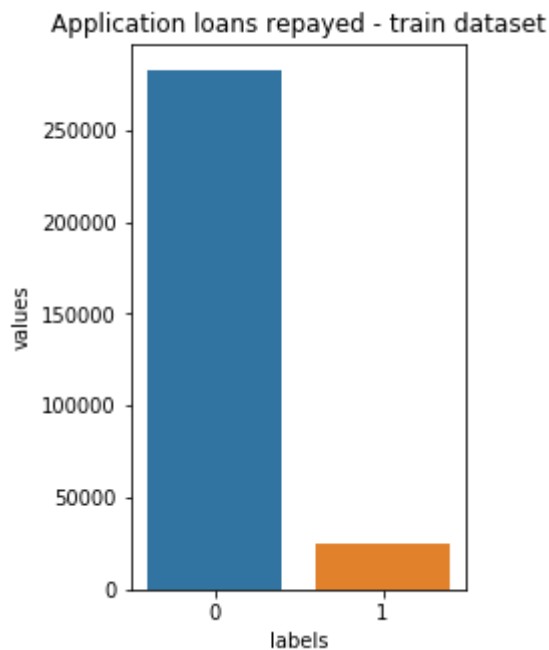
As we can see, the dataset is imbalanced. Positive class ratio is only 0.081. We can also see the number of customers in each label.

In [6]:

```

1 temp = train_set["TARGET"].value_counts()
2 df = pd.DataFrame({'labels': temp.index,
3                   'values': temp.values
4                   })
5 plt.figure(figsize = (3,5))
6 plt.title('Application loans repayed - train dataset')
7 sns.set_color_codes("pastel")
8 sns.barplot(x = 'labels', y="values", data=df)
9 locs, labels = plt.xticks()
10 plt.show()

```



Since it is an imbalanced we should check some of the features distributions for each class. Distributions or ratios might be similar or might show different behaviors for each class.

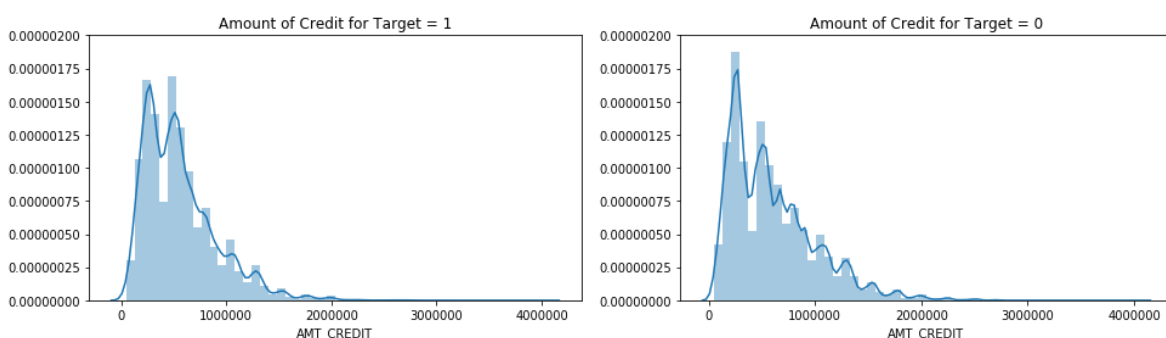
## Distribution of credit amount of the loan for positive and negative classes

In [11]:

```

1 fig, axs = plt.subplots(1,2, figsize=(16, 4))
2 axs[0].set_ylim([0, 0.000002])
3 axs[1].set_ylim([0, 0.000002])
4 ax = sns.distplot(train_set_Target1["AMT_CREDIT"].dropna(),ax=axs[0]).set_title("Amount of Credit for Target = 1")
5 ay = sns.distplot(train_set_Target0["AMT_CREDIT"].dropna(),ax=axs[1]).set_title("Amount of Credit for Target = 0")

```

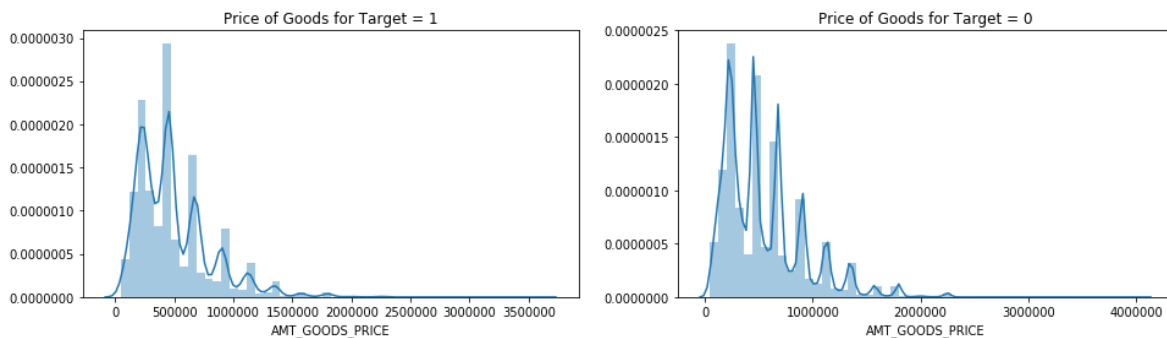


Credit amount of the loan shows similar behaviors for both classes, both are right skewed. Most of the credit amounts are less than 1,000,000. The average amount of credit is higher for negative class (mean = 602648) than the positive class (mean = 557779).

## Distribution of goods price for given loan for positive and negative classes

In [21]:

```
1 fig, axs = plt.subplots(1,2, figsize=(16, 4))
2
3 ax = sns.distplot(train_set_Target1["AMT_GOODS_PRICE"].dropna(),ax=axs[0]).set_title("Price of Goods for Target = 1")
4 ay = sns.distplot(train_set_Target0["AMT_GOODS_PRICE"].dropna(),ax=axs[1]).set_title("Price of Goods for Target = 0")
```

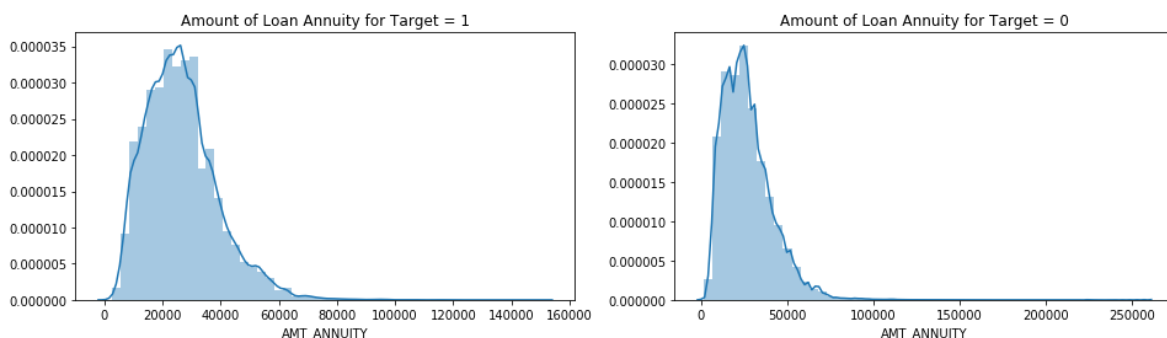


Price of the goods for which the loan is given shows similar behaviors, both are right skewed. Most of the price of the goods are less than 1,500,000. The average price of the goods credit is higher for negative class (mean = 542737) than the positive class (mean = 488972).

## Distribution of loan annuity for positive and negative classes

In [23]:

```
1 fig, axs = plt.subplots(1,2, figsize=(16, 4))
2
3 ax = sns.distplot(train_set_Target1["AMT_ANNUIITY"].dropna(),ax=axs[0]).set_title("Amount of Loan Annuity for Target = 1")
4 ay = sns.distplot(train_set_Target0["AMT_ANNUIITY"].dropna(),ax=axs[1]).set_title("Amount of Loan Annuity for Target = 0")
```



We can see that amount of annuity for the loan also shows similar behaviors for positive and negative classes.

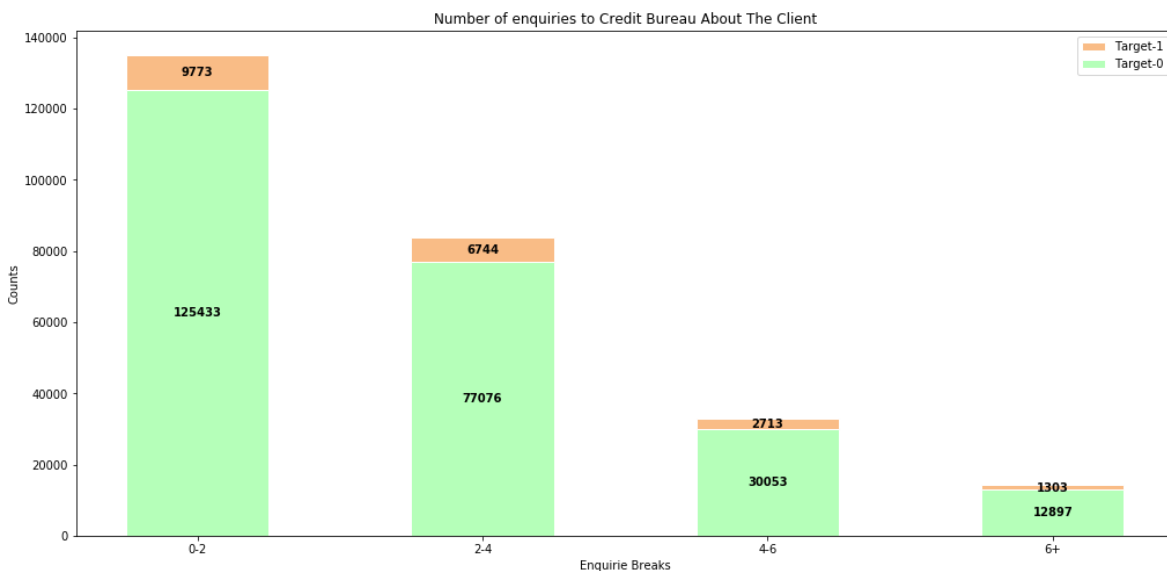
# Number of enquiries to Credit Bureau about the client in a year

In [26]:

```

1 Df_count = train_set.groupby(['TARGET', 'AMT_REQ_CREDIT_BUREAU_YEAR']).size().reset_index()
2 train_set_0 = Df_count[Df_count['TARGET'] == 0]
3 train_set_1 = Df_count[Df_count['TARGET'] == 1]
4
5 list_0 = []
6 list_1 = []
7 for i in [2,4,6,50]:
8     c0 = train_set_0.loc[(train_set_0['AMT_REQ_CREDIT_BUREAU_YEAR'] < i), 'counts'].sum()
9     c1 = train_set_1.loc[(train_set_1['AMT_REQ_CREDIT_BUREAU_YEAR'] < i), 'counts'].sum()
10    list_0.append(c0)
11    list_1.append(c1)
12 list_0[1:len(list_0)] = [j-i for i, j in zip(list_0[:-1], list_0[1:])]
13 list_1[1:len(list_1)] = [j-i for i, j in zip(list_1[:-1], list_1[1:])]
14
15 plt.figure(figsize=(17,8))
16 r = [0,1,2,3]
17 barWidth = 0.5
18 names = ('0-2', '2-4', '4-6', '6+')
19 # Create green Bars
20 p0 = plt.bar(r, list_0, color='#b5ffb9', edgecolor='white', width=barWidth)
21 p1 = plt.bar(r, list_1, bottom=list_0, color='#f9bc86', edgecolor='white', width=barWidth)
22 plt.xticks(r, names)
23 plt.legend((p1[0], p0[0]), ('Target-1', 'Target-0'))
24 plt.title('Number of enquiries to Credit Bureau About The Client')
25 plt.xlabel('Enquirie Breaks')
26 plt.ylabel('Counts')
27
28 for r1, r2 in zip(p0, p1):
29     h1 = r1.get_height()
30     h2 = r2.get_height()
31     plt.text(r1.get_x() + r1.get_width() / 2., h1 / 2., "%d" % h1, ha="center", va="center")
32     plt.text(r2.get_x() + r2.get_width() / 2., h1 + h2 / 2., "%d" % h2, ha="center", va="center")
33
34 plt.show()

```



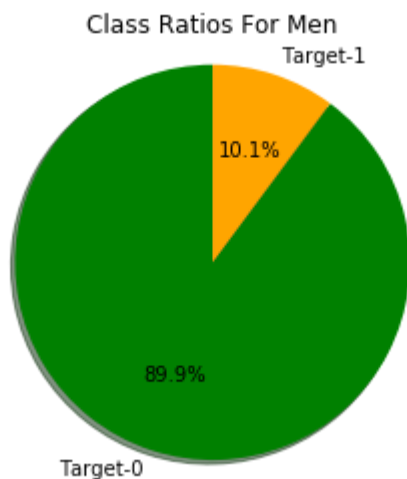


There is a large number of customers having enquiries to Credit Bureau between 0-2. As the number of enquiries to Credit Bureau increases, number of occurrence decreases in both classes. Number of enquiries is much greater for Target-0 than Target-1.

## Positive and Negative class ratios for men and women

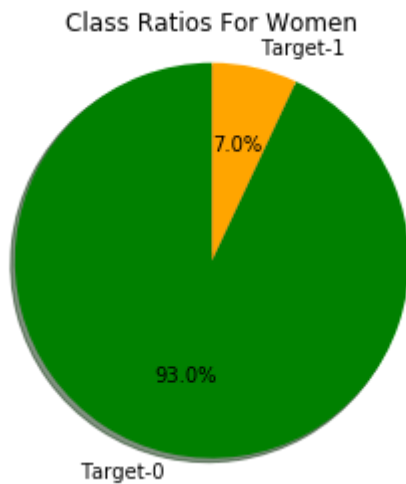
In [35]:

```
1 CashLoans = train_set[train_set.CODE_GENDER == 'M']
2 temp = CashLoans["TARGET"].value_counts()
3 df = pd.DataFrame({'labels': ["Target-0", "Target-1"],
4                       'values': temp.values
5                     })
6 colors = ['green', 'orange']
7 plt.pie(df["values"], labels=df["labels"], colors=colors,
8         autopct='%1.1f%%', shadow=True, startangle=90)
9 plt.title("Class Ratios For Men")
10
11 plt.axis('equal')
12 plt.show()
```



In [36]:

```
1 CashLoans = train_set[train_set.CODE_GENDER == 'F']
2 temp = CashLoans["TARGET"].value_counts()
3 df = pd.DataFrame({'labels': ["Target-0", "Target-1"],
4                       'values': temp.values
5                       })
6 colors = ['green', 'orange']
7 plt.pie(df["values"], labels=df["labels"], colors=colors,
8         autopct='%1.1f%%', shadow=True, startangle=90)
9 plt.title("Class Ratios For Women")
10
11 plt.axis('equal')
12 plt.show()
```



Positive class ratio is higher among men than women.

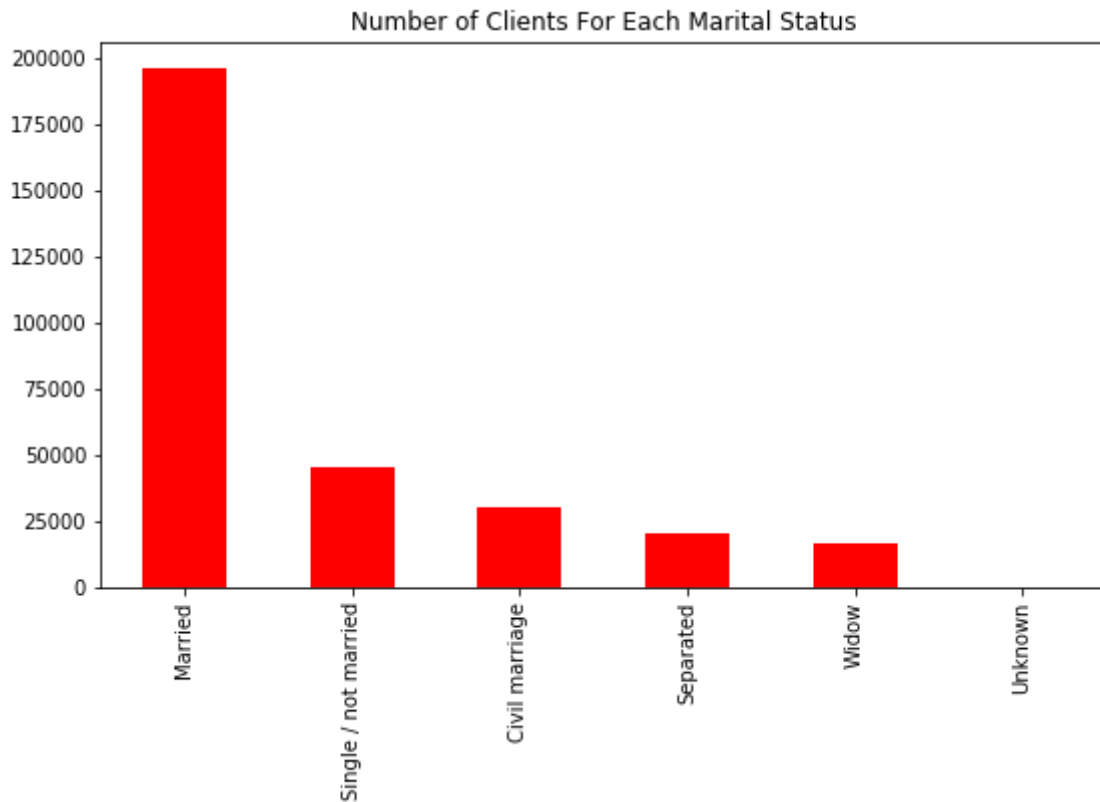
## Categorical Distributions :

1. Family Status
2. Contract Type
3. Education
4. Occupation

## Distribution of Family Status

In [16]:

```
1 fms = train_set['NAME_FAMILY_STATUS'].value_counts().plot(kind='bar',  
2                               figsize=(9,5),  
3                               title="Number of Clients For Each Marital Status",
```

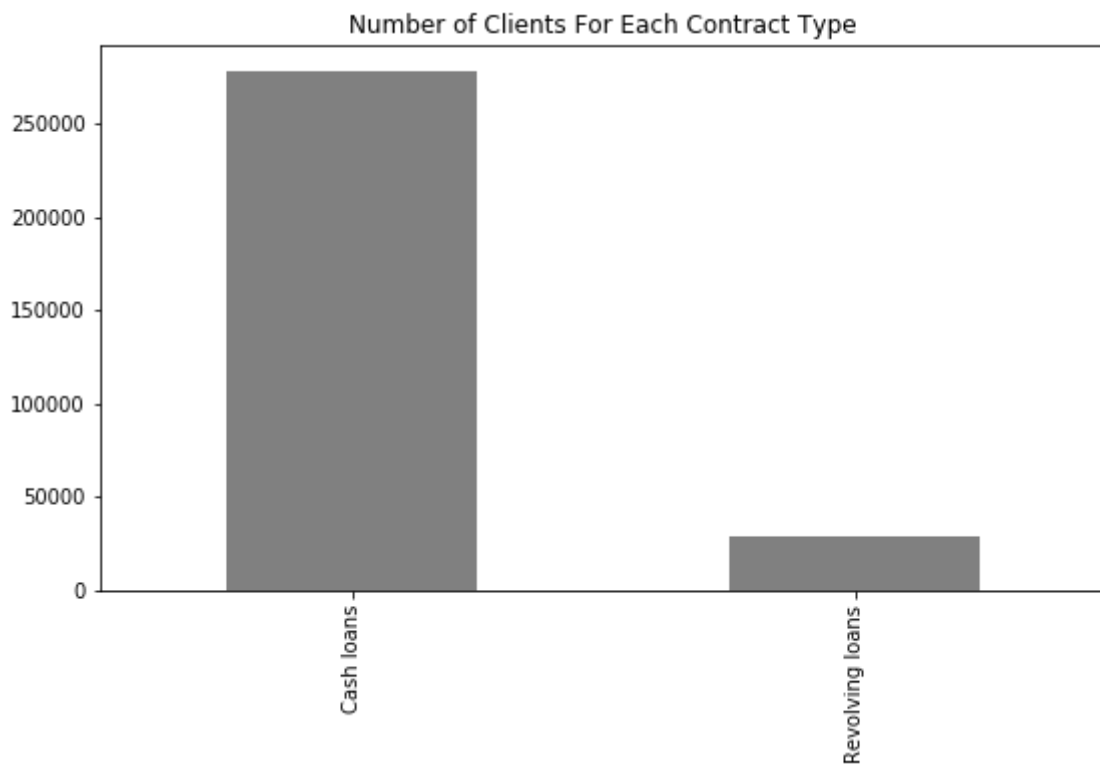


Mostly married people have applied for a larger number of loan applications around 200K, followed by Single/not married and civil marriage.

## Distribution of Contract Type

In [4]:

```
1 loantyp = train_set['NAME_CONTRACT_TYPE'].value_counts().plot(kind='bar',  
2                                     figsize=(9,5),  
3                                     title="Number of Clients For Each Contract Type",
```

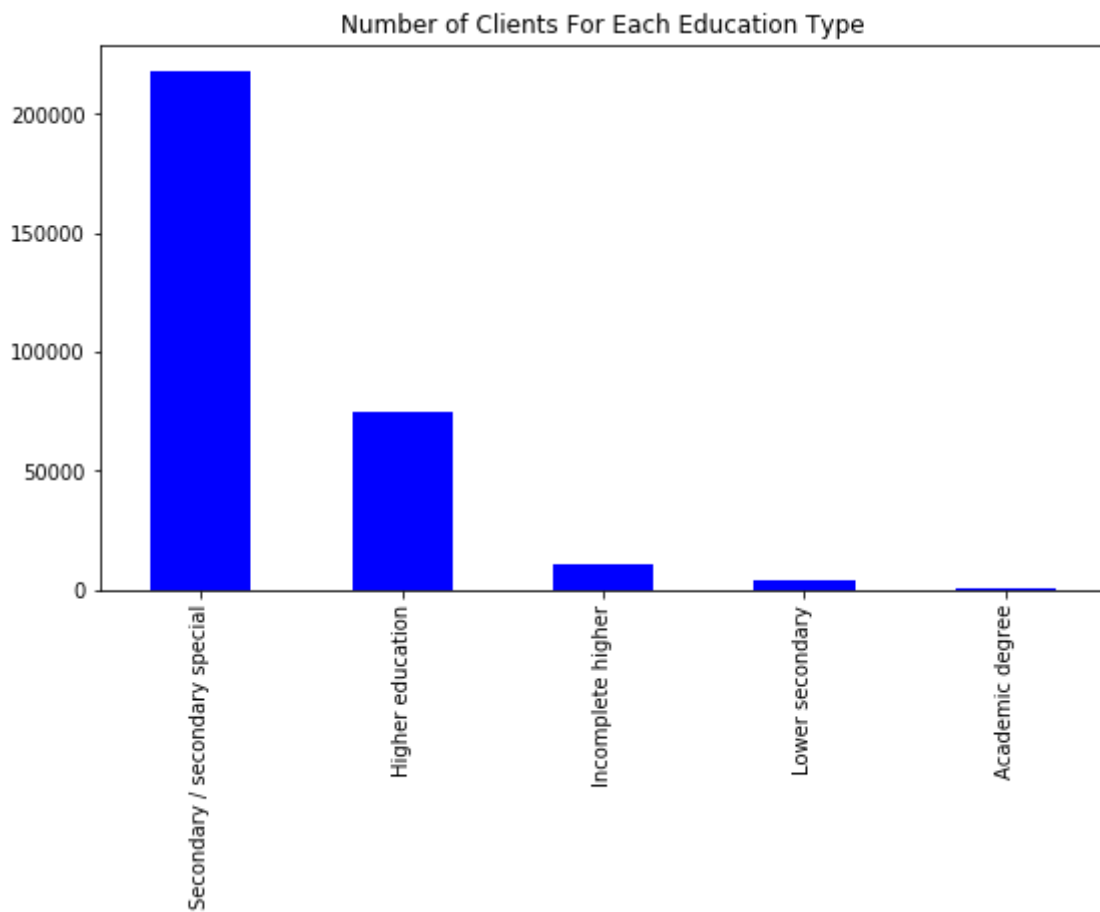


Majority of the clients are applying for cash loan.

## Distribution of Education Type

In [24]:

```
1 edctyp = train_set['NAME_EDUCATION_TYPE'].value_counts().plot(kind='bar',  
2                                     figsize=(9,5),  
3                                     title="Number of Clients For Each Education Type",
```



The number of people applying for a loan among highly educated people is lower. Majority of applicants belong to secondary education type.

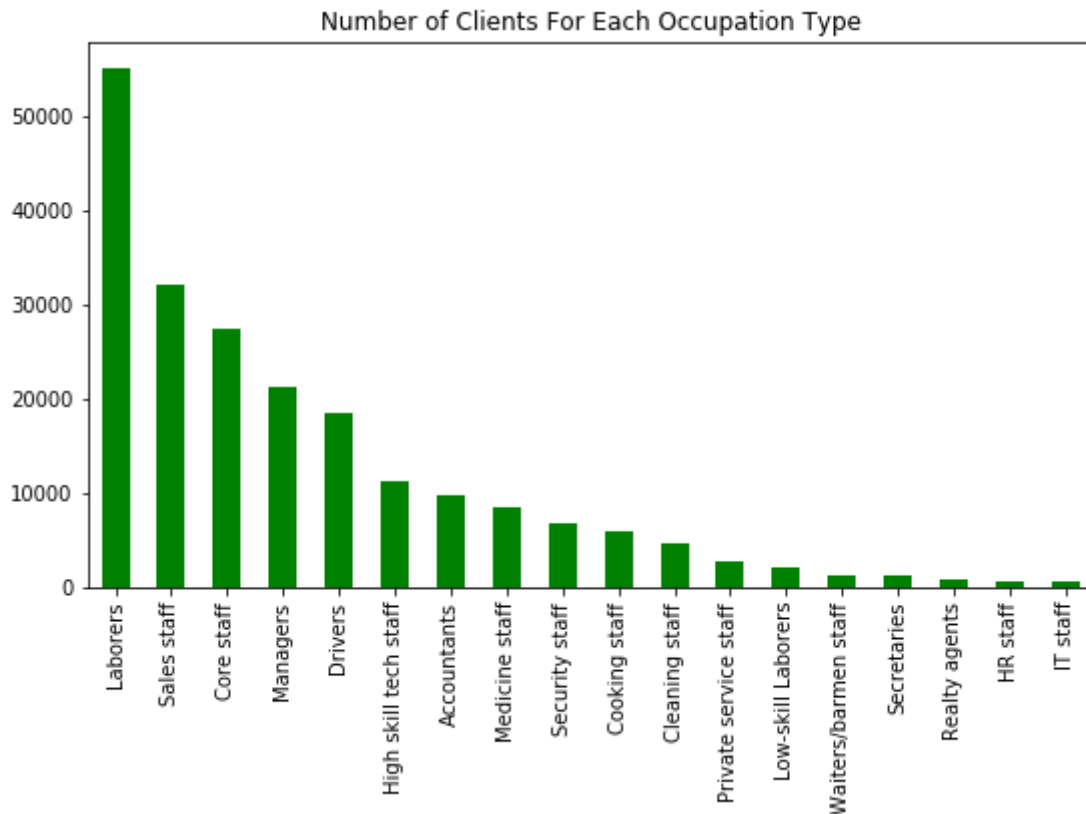
## Distribution of Occupation Type

In [5]:

```

1  ocptyp = train_set['OCCUPATION_TYPE'].value_counts().plot(kind='bar',
2                                     figsize=(9,5),
3                                     title="Number of Clients For Each Occupation Type")

```



Most of the loans are taken by Laborers, followed by Sales staff and core staff. IT staff take the lowest amount of loans.

## 2. Feature Generation

### Feature Generation For Train Set

There is one main data source (application\_train.csv), including information about each loan application with a unique loan id. It has information about the clients such as clients' gender, if s/he owns a house, number of children the client has, income of the client, scores from each client from external sources and, credit amount and loan annuity of the loan that the client applies for. From the provided columns, we can extract other columns such as; a flag showing if a customer has greater income than the credit he applies, the ratio of the credit amount over the income amount of the client, the ratio of the annuity amount over the income amount of the client, the ratio of the credit amount over the annuity amount, the ratio of the days s/he's employed in his life, and with which ratio of the amount of the goods price, he uses credit. We can add these features to the main dataset. The code to attach these features is shown below.

In [24]:

```
1 #Feature Generation From Application Data (main data source):
2 #If a customer has greater income than the credit he applies
3 train_set['INCOME_GT_CREDIT_FLAG'] = np.where(train_set['AMT_INCOME_TOTAL'] > train_set['AMT_CREDIT'], 1, 0)
4
5 #Credit Income Ratio
6 train_set['CREDIT_INCOME_PERCENT'] = train_set['AMT_CREDIT'] / train_set['AMT_INCOME_TOTAL']
7
8 #Annuity Income Ratio
9 train_set['ANNUITY_INCOME_PERCENT'] = train_set['AMT_ANNUITY'] / train_set['AMT_INCOME_TOTAL']
10
11 # Column to represent Credit Term
12 train_set['CREDIT_TERM'] = train_set['AMT_CREDIT'] / train_set['AMT_ANNUITY']
13
14 # Column to represent Days Employed percent in his life
15 train_set['DAYS_EMPLOYED_PERCENT'] = train_set['DAYS_EMPLOYED'] / train_set['DAYS_BIRTH']
16
17 #Credit-Price of Goods Ratio
18 train_set['CREDIT_PRICE_OF_GOODS_RATIO'] = train_set['AMT_CREDIT'] / train_set['AMT_GOODS_PRICE']
```

There are other supplementary datasets that we can combine to the main dataset and generate other features. Each supplementary dataset has information about the unique loan id, but multiple instances for a single loan id. Once we group the data by the unique loan id and merge it to the train set; we can have information about the client's previous behaviors and other actions related to the loan. The code of how to use each supplementary data source and what features to add from it is shown below with illustrative comments.

In [25]:

```

1  #Feature Generation From Other Data (other data sources):
2  ##-----
3  #-----Combining Breau Data to Train set-----
4  #group data by id: We can extract number of previous credit, number of active credit,
5  bureau['CREDIT_ACTIVE_FLAG'] = np.where(bureau['CREDIT_ACTIVE'] == 'Active', 1, 0) #ta
6  BreauSumTable = bureau[['SK_ID_CURR', 'CREDIT_ACTIVE_FLAG', 'AMT_CREDIT_SUM', 'AMT_CREDI
7  BreauSumTable = BreauSumTable.rename(columns={'CREDIT_ACTIVE_FLAG': 'NUMBER_OF_ACTIVE_C
8
9  BreauCountTable = bureau[['SK_ID_CURR', 'SK_ID_BUREAU']].groupby(['SK_ID_CURR']).count(
10 BreauCountTable = BreauCountTable.rename(columns={'SK_ID_BUREAU': 'NUMBER_OF_PREVIOUS_C
11
12 #merge grouped data into train set:
13 train_set = train_set.merge(BreauSumTable, on='SK_ID_CURR', how='left')
14 train_set = train_set.merge(BreauCountTable, on='SK_ID_CURR', how='left')
15
16 #fill ne values with zero
17 train_set.update(train_set[BreauSumTable.columns].fillna(0))
18 train_set.update(train_set[BreauCountTable.columns].fillna(0))
19 #train_set.head()
20
21
22 ##-----
23 #-----Combining pos_cash Data to Train set-----
24 #group data by id: We can extract number of previous credit, number of active credit,
25 pos_cash['ACTIVE_STATUS'] = np.where(pos_cash['NAME_CONTRACT_STATUS'] == 'Active', 1,
26 pos_cashSum = pos_cash[['SK_ID_CURR', 'ACTIVE_STATUS']].groupby(['SK_ID_CURR']).sum().r
27 pos_cashSum = pos_cashSum.rename(columns={'ACTIVE_STATUS': 'NUMBER_OF_ACTIVE_CONTRACT_C
28
29 pos_cashCount = pos_cash[['SK_ID_CURR', 'SK_ID_PREV']].groupby(['SK_ID_CURR']).count().
30 pos_cashCount = pos_cashCount.rename(columns={'SK_ID_PREV': 'NUMBER_OF_PREVIOUS_CREDIT_
31
32 pos_cashAvg = pos_cash[['SK_ID_CURR', 'CNT_INSTALLMENT', 'CNT_INSTALLMENT_FUTURE']].groupb
33 pos_cashAvg = pos_cashAvg .rename(columns={'CNT_INSTALLMENT': 'AVG_OF_INSTALLMENTS_CASH'
34
35 #merge grouped data into train set:
36 train_set = train_set.merge(pos_cashCount, on='SK_ID_CURR', how='left')
37 train_set = train_set.merge(pos_cashSum, on='SK_ID_CURR', how='left')
38 train_set = train_set.merge(pos_cashAvg, on='SK_ID_CURR', how='left')
39
40 #fill ne values with zero
41 train_set.update(train_set[pos_cashCount.columns].fillna(0))
42 train_set.update(train_set[pos_cashSum.columns].fillna(0))
43 train_set.update(train_set[pos_cashAvg.columns].fillna(0))
44 #train_set.head()
45
46
47 ##-----
48 #-----Combining credit_card_balance Data to Train set-----
49
50 #group data by id: We can extract number of avtive card contracts, number of previous
51 #//total balance, total limit, total drawings amounts and count of drawings of each cu
52 credit_card_balance['ACTIVE_STATUS'] = np.where(credit_card_balance['NAME_CONTRACT_STA
53 credit_card_balanceSum = credit_card_balance[['SK_ID_CURR', 'ACTIVE_STATUS']].groupby([
54 credit_card_balanceSum = credit_card_balanceSum.rename(columns={'ACTIVE_STATUS': 'NUMBE
55
56 credit_card_balanceCount = credit_card_balance[['SK_ID_CURR', 'SK_ID_PREV']].groupby(['
57 credit_card_balanceCount = credit_card_balanceCount.rename(columns={'SK_ID_PREV': 'NUMB
58
59 credit_card_balanceTotal = credit_card_balance[['SK_ID_CURR', 'AMT_BALANCE', 'AMT_CREDIT

```



```

60 credit_card_balanceTotal = credit_card_balanceTotal .rename(columns={'AMT_BALANCE': 'TO
61
62 #merge groupped data into train set:
63 train_set = train_set.merge(credit_card_balanceSum, on='SK_ID_CURR', how='left')
64 train_set = train_set.merge(credit_card_balanceCount, on='SK_ID_CURR', how='left')
65 train_set = train_set.merge(credit_card_balanceTotal, on='SK_ID_CURR', how='left')
66
67 #fill ne values with zero
68 train_set.update(train_set[pos_cashCount.columns].fillna(0))
69 train_set.update(train_set[pos_cashSum.columns].fillna(0))
70 train_set.update(train_set[pos_cashAvg.columns].fillna(0))
71 #train_set.head()
72
73
74 ##-----
75 #-----Combining previous_applicaton Data to Train set
76
77 #group data by id: We can extract count of previous application, count of approved apl
78 #//total annuity, total application, total credit of each customer
79 previous_applicaton['ACTIVE_STATUS'] = np.where(previous_applicaton['NAME_CONTRACT_STA
80 previous_applicatonSum = previous_applicaton[['SK_ID_CURR', 'ACTIVE_STATUS']].groupby([
81 previous_applicatonSum = previous_applicatonSum.rename(columns={'ACTIVE_STATUS': 'COUNT
82
83 previous_applicatonCount = previous_applicaton[['SK_ID_CURR', 'SK_ID_PREV']].groupby(['
84 previous_applicatonCount = previous_applicatonCount.rename(columns={'SK_ID_PREV': 'NUMB
85
86 previous_applicatonTotal = previous_applicaton[['SK_ID_CURR', 'AMT_ANNUITY', 'AMT_APPLI
87 previous_applicatonTotal = previous_applicatonTotal .rename(columns={'AMT_ANNUITY': 'TO
88
89 #merge groupped data into train set:
90 train_set = train_set.merge(previous_applicatonSum, on='SK_ID_CURR', how='left')
91 train_set = train_set.merge(previous_applicatonCount, on='SK_ID_CURR', how='left')
92 train_set = train_set.merge(previous_applicatonTotal, on='SK_ID_CURR', how='left')
93
94 #fill ne values with zero
95 train_set.update(train_set[previous_applicatonCount.columns].fillna(0))
96 train_set.update(train_set[previous_applicatonSum.columns].fillna(0))
97 train_set.update(train_set[previous_applicatonTotal.columns].fillna(0))
98 #train_set.head()
99
100
101 ##-----
102 #-----Combining installments_payments Data to Train set
103
104 #group data by id: We can extract average number of installments that each customer ma
105 installments_paymentsAvg = installments_payments[['SK_ID_CURR', 'NUM_INSTALLMENT_NUMBER'
106 installments_paymentsAvg = installments_paymentsAvg.rename(columns={'NUM_INSTALLMENT_NU
107
108 #merge groupped data into train set:
109 train_set = train_set.merge(installments_paymentsAvg, on='SK_ID_CURR', how='left')
110
111 #fill ne values with zero
112 train_set.update(train_set[installments_paymentsAvg.columns].fillna(0))

```

Some columns have a considerable amount of missing values which can be misleading for the predictions of the model on the test set. Due to this reason, I decided to remove columns having more than 10% of null values in it. After composing our final train data set by adding and removing features, we can write this data frame into a csv file for further use.

In [26]:

```

1  #-----Feature Elimination:
2
3  #Removing Columns Having Null Valio of ratio more than 0.1
4  def rmissingvaluecol(dff,threshold):
5      l = []
6      l = list(dff.drop(dff.loc[:,list((100*(dff.isnull().sum()/len(dff.index))>=threshold))].index))
7      print("# Columns having more than %s percent missing values:"%threshold,(dff.shape[1]-len(l)))
8      print("Columns:\n",list(set(list((dff.columns.values))) - set(l)))
9      return l
10
11
12  remaningColumns= rmissingvaluecol(train_set,10) #Here threshold is 10% which means we are removing columns having more than 10% missing values
13  train_set2 = train_set[remaningColumns]
14
15  train_set2.to_csv(directory+ '\\train_set_final.csv')

```

# Columns having more than 10 percent missing values: 63

Columns:

```

['COMMONAREA_MEDI', 'YEARS_BEGINEXPLUATATION_AVG', 'YEARS_BEGINEXPLUATATION_MODE', 'HOUSETYPE_MODE', 'OCCUPATION_TYPE', 'FLOORSMAX_MODE', 'ELEVATORS_AVG', 'AMT_REQ_CREDIT_BUREAU_MON', 'AMT_REQ_CREDIT_BUREAU_QRT', 'COUNT_OF_DRAWINGS_CARD', 'LIVINGAPARTMENTS_MODE', 'APARTMENTS_AVG', 'TOTALAREA_MODE', 'OWN_CAR_AGE', 'AMT_REQ_CREDIT_BUREAU_YEAR', 'FLOORSMIN_MEDI', 'COMMONAREA_AVG', 'NUMBER_OF_ACTIVE_CONTRACT_CARD', 'YEARS_BUILD_AVG', 'FLOORSMIN_AVG', 'AMT_REQ_CREDIT_BUREAU_WEEK', 'COMMONAREA_MODE', 'YEARS_BUILD_MEDI', 'FLOORSMAX_AVG', 'LIVINGAPARTMENTS_AVG', 'ELEVATORS_MODE', 'LIVINGAREA_AVG', 'NONLIVINGAPARTMENTS_MEDI', 'BASEMENTAREA_MODE', 'ELEVATORS_MEDI', 'YEARS_BEGINEXPLUATATION_MEDI', 'NONLIVINGAPARTMENTS_MODE', 'ENTRANCES_AVG', 'FLOORSMAX_MEDI', 'NONLIVINGAREA_MEDI', 'EXT_SOURCE_1', 'NUMBER_OF_PREVIOUS_CARD', 'AMT_REQ_CREDIT_BUREAU_DAY', 'NONLIVINGAREA_AVG', 'BASEMENTAREA_AVG', 'LIVINGAREA_MODE', 'WALLSMATERIAL_MODE', 'YEARS_BUILD_MODE', 'NONLIVINGAREA_MODE', 'AMT_REQ_CREDIT_BUREAU_HOUR', 'EXT_SOURCE_3', 'TOTAL_LIMIT_CARD', 'TOTAL_AMOUNT_BALANCE_CARD', 'LANDAREA_AVG', 'APARTMENTS_MEDI', 'APARTMENTS_MODE', 'LANDAREA_MODE', 'FLOORSMIN_MODE', 'FONDKAPREMONT_MODE', 'LIVINGAREA_MEDI', 'ENTRANCES_MEDI', 'ENTRANCES_MODE', 'BASEMENTAREA_MEDI', 'LANDAREA_MEDI', 'NONLIVINGAPARTMENTS_AVG', 'TOTAL_DRAWINGS_CARD', 'LIVINGAPARTMENTS_MEDI', 'EMERGENCYSTATE_MODE']

```

## Feature Generation For Test Set

We should follow the same steps as we did for train set, on the test set as well, so we can use the test set for making predictions. Test set does not have the TARGET feature.

In [27]:

```

1  ###-----Test Set Preparation (with the columns from test set):
2
3  test_set['INCOME_GT_CREDIT_FLAG'] = np.where(test_set['AMT_INCOME_TOTAL'] > test_set['INCOME_TOTAL'], 1, 0)
4
5  #Credit Income Ratio
6  test_set['CREDIT_INCOME_PERCENT'] = test_set['AMT_CREDIT'] / test_set['AMT_INCOME_TOTAL']
7
8  #Annuity Income Ratio
9  test_set['ANNUITY_INCOME_PERCENT'] = test_set['AMT_ANNUITY'] / test_set['AMT_INCOME_TOTAL']
10
11 # Column to represent Credit Term
12 test_set['CREDIT_TERM'] = test_set['AMT_CREDIT'] / test_set['AMT_ANNUITY']
13
14 # Column to represent Days Employed percent in his life
15 test_set['DAYS_EMPLOYED_PERCENT'] = test_set['DAYS_EMPLOYED'] / test_set['DAYS_BIRTH']
16
17 #Credit-Price of Goods Ratio
18 test_set['CREDIT_PRICE_OF_GOODS_RATIO'] = test_set['AMT_CREDIT'] / test_set['AMT_GOODS_PRICE']
19
20 ##-----
21 #merge grouped data into train set:
22 test_set = test_set.merge(BreauSumTable, on='SK_ID_CURR', how='left')
23 test_set = test_set.merge(BreauCountTable, on='SK_ID_CURR', how='left')
24
25 #fill ne values with zero
26 test_set.update(test_set[BreauSumTable.columns].fillna(0))
27 test_set.update(test_set[BreauCountTable.columns].fillna(0))
28 #test_set.head()
29
30 #merge grouped data into train set:
31 test_set = test_set.merge(pos_cashCount, on='SK_ID_CURR', how='left')
32 test_set = test_set.merge(pos_cashSum, on='SK_ID_CURR', how='left')
33 test_set = test_set.merge(pos_cashAvg, on='SK_ID_CURR', how='left')
34
35 #fill ne values with zero
36 test_set.update(test_set[pos_cashCount.columns].fillna(0))
37 test_set.update(test_set[pos_cashSum.columns].fillna(0))
38 test_set.update(test_set[pos_cashAvg.columns].fillna(0))
39 #test_set.head()
40
41 #merge grouped data into train set:
42 test_set = test_set.merge(credit_card_balanceSum, on='SK_ID_CURR', how='left')
43 test_set = test_set.merge(credit_card_balanceCount, on='SK_ID_CURR', how='left')
44 test_set = test_set.merge(credit_card_balanceTotal, on='SK_ID_CURR', how='left')
45
46 #fill ne values with zero
47 test_set.update(test_set[pos_cashCount.columns].fillna(0))
48 test_set.update(test_set[pos_cashSum.columns].fillna(0))
49 test_set.update(test_set[pos_cashAvg.columns].fillna(0))
50 #test_set.head()
51
52 #merge grouped data into train set:
53 test_set = test_set.merge(previous_applicatonSum, on='SK_ID_CURR', how='left')
54 test_set = test_set.merge(previous_applicatonCount, on='SK_ID_CURR', how='left')
55 test_set = test_set.merge(previous_applicatonTotal, on='SK_ID_CURR', how='left')
56
57 #fill ne values with zero
58 test_set.update(test_set[previous_applicatonCount.columns].fillna(0))
59 test_set.update(test_set[previous_applicatonSum.columns].fillna(0))

```

```
60 test_set.update(test_set[previous_applicatonTotal.columns].fillna(0))
61 #test_set.head()
62
63 #merge groupped data into train set:
64 test_set = test_set.merge(installments_paymentsAvg, on='SK_ID_CURR', how='left')
65
66 #fill ne values with zero
67 test_set.update(test_set[installments_paymentsAvg.columns].fillna(0))
68
69
70 #Removes TARGET from remaining columes
71 del remainingColumns[1]
72 test_set2 = test_set[remainingColumns]
73 test_set2.to_csv(directory+ '\\test_set_final.csv')
```

### 3. Model Application

Home Credit Default Risk Project is a binary classificaiton problem. However making binary predictions for the target will not be meaningful. We should use probabilities instead, to score the clients according to how capable they are to repay the loan. For example, we should not classify two clients having scores 0.75 and 0.9 as positive just because they both have greater score than 0.5. We should consider the client with score 0.9 has higher probability on positive class than the client with score 0.75.

Since it is a long-term predictions problem, the algorithm can take some time to run. That is the reason why we can use ensemble methods.

We can use Area under Curve to summarize the model performance, as it is the most expressive way. We can plot the ROC curve, to visualize the performance of a binary classifier.

### Validation of the model using only train set and splitting it into two sets:

We have a train set with 307511 rows and 86 features. As train set and test sets for this project are split on the loan id, we can split the train set into smaller train and test sets for the application of the model and to see the model preformance.

In this section, I split train data into two. I prepared both sets for the catboost algorithm. I fitted the model on the 80% of the data, and predict the targets for the remaining 20%. I used 5 folds for cross validaiton and made 500 itearions with early stopping 25. I used AUC as a measurement of the model performance and plotted ROC curve to visualize.

In [35]:

```

1 import pandas as pd
2 import catboost as cat
3 from sklearn.preprocessing import LabelEncoder
4 import os
5 import pickle
6 from sklearn.metrics import roc_auc_score
7 from sklearn.model_selection import train_test_split
8 from sklearn import metrics
9 import matplotlib.pyplot as plt
10
11
12 directory = 'C:\\CreditRiskProject\\Datasets'
13 train_set = pd.read_csv(directory+ '\\train_set_final.csv', index_col=0)

```

In [38]:

```
1 train_set.shape
```

Out[38]:

(307511, 86)

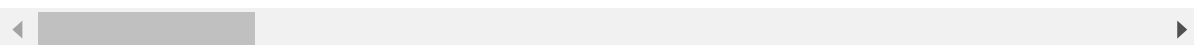
In [39]:

```
1 train_set.head()
```

Out[39]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY
0	100002	1	Cash loans	M	N	N
1	100003	0	Cash loans	F	N	N
2	100004	0	Revolving loans	M	Y	N
3	100006	0	Cash loans	F	N	N
4	100007	0	Cash loans	M	N	N

5 rows × 86 columns



As we can see from the table above, in this data set, there are categorical values that we should encode to fit the model. Then we can apply CatBoost algorithm and measure the performance by AUC and see the feature importances. The code is provided below with illustrative comments.

In [40]:

```

1  #Remove Unnecessary columns
2  def removenonuniquecol(dataset):
3      dropcols = [col for col in dataset.columns if dataset[col].nunique(dropna=True)==1]
4      print ('Removing columns: ',dropcols)
5      dataset.drop(dropcols,axis= 1,inplace= True,errors= 'ignore')
6      return dropcols
7
8  #For Label encoding
9  def labelencoder(dataset):
10     objectlist = dataset.select_dtypes(include=['object']).copy()
11     cat_col = [col for col in dataset.columns if col in objectlist]
12     for col in cat_col:
13         print("Encoding ",col)
14         lbl = LabelEncoder()
15         dataset[col].fillna(-999)
16         lbl.fit(list(dataset[col].values.astype('str')))
17         dataset[col] = lbl.transform(list(dataset[col].values.astype('str')))
18     return cat_col
19
20 #####Trainset preperation#####
21
22 def TrainPrep(Datasetname):
23     removedCols = removenonuniquecol(Datasetname)
24     Predictors = [col for col in Datasetname]
25     Predictors = [col for col in Predictors if col not in removedCols]
26     Predictors = [col for col in Predictors if col not in ['TARGET']]
27     DfLabel = Datasetname['TARGET']
28     encodedList = labelencoder(Datasetname)
29     return Predictors,Datasetname, DfLabel, removedCols, encodedList
30
31 Predictors,Df_1, Label_1, removedCols, encodedList = TrainPrep(train_set)
32
33 X_train, X_test, y_train, y_test = train_test_split(Df_1, Label_1, stratify=Df_1['TARGET'])
34
35
36 #####Training Classifier#####
37
38
39 def catboosttrainer(X,y,features,initparam,modelname,modelpath,docpath,cvfold = 5):
40     print ("searching for optimal iteration count...")
41     trainpool = cat.Pool(X[features],y)
42     cvresult = cat.cv(params= initparam, fold_count=cvfold, pool=trainpool,stratified = True)
43     initparam['iterations'] = (len(cvresult)) - (initparam['od_wait']+1)
44     del initparam['od_wait']
45     del initparam['od_type']
46     print ("optimal iteration count is ", initparam['iterations'])
47     print ("fitting model...")
48     clf = cat.CatBoostClassifier(** initparam)
49     clf.fit(trainpool)
50     imp = clf.get_feature_importance(trainpool,fstr_type='FeatureImportance')
51     dfimp = pd.DataFrame(imp,columns = ['CatBoostImportance'])
52     dfimp.insert(0,column='Feature', value=features)
53     dfimp = dfimp.sort_values(['CatBoostImportance','Feature'], ascending= False)
54     xlsxpath = os.path.join(docpath,modelname+".xlsx")
55     dfimp.to_excel(xlsxpath)
56     print ("pickling model...")
57     picklepath = os.path.join(modelpath,modelname)
58     with open(picklepath,'wb') as fout:
59         pickle.dump(clf, fout)

```

```

60     return cvresult,clf,initparam,dfimp
61
62
63
64 modelpath = 'C:\\CreditRiskProject'
65 docpath = 'C:\\CreditRiskProject'
66 CatBoostParam = { 'iterations': 500, 'od_type': 'Iter', 'od_wait': 25, 'loss_function':
67
68 cvresult,clf,initparam,dfimp = catboosttrainer(X_train,y_train,Predictors,CatBoostParam
69
70 proba = clf.predict_proba(X_test[Predictors])[:,1]
71 auc = roc_auc_score(y_test,proba)
72 print(auc)
73
74

```

```

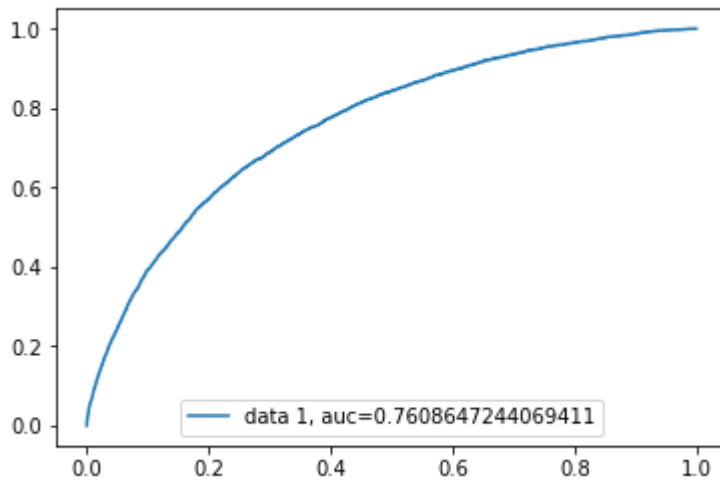
Removing columns: []
Encoding  NAME_CONTRACT_TYPE
Encoding  CODE_GENDER
Encoding  FLAG_OWN_CAR
Encoding  FLAG_OWN_REALTY
Encoding  NAME_TYPE_SUITE
Encoding  NAME_INCOME_TYPE
Encoding  NAME_EDUCATION_TYPE
Encoding  NAME_FAMILY_STATUS
Encoding  NAME_HOUSING_TYPE
Encoding  WEEKDAY_APPR_PROCESS_START
Encoding  ORGANIZATION_TYPE
searching for optimal iteration count...
0:      test: 0.6002952 best: 0.6002952 (0)
1:      test: 0.6478499 best: 0.6478499 (1)
2:      test: 0.6599599 best: 0.6599599 (2)
3:      test: 0.6717097 best: 0.6717097 (3)
4:      test: 0.6805417 best: 0.6805417 (4)
5:      test: 0.6882623 best: 0.6882623 (5)
6:      test: 0.6947376 best: 0.6882623 (6)

```

Feature importances are provided in an excel file (CBmodel.xlsx). AUC of the model is 76.1%. You can see the ROC curve provided below.

In [42]:

```
1 fpr, tpr, _ = metrics.roc_curve(y_test, proba)
2 auc = metrics.roc_auc_score(y_test, proba)
3 plt.plot(fpr,tpr,label="data 1, auc="+str(auc))
4 plt.legend(loc=8)
5 plt.show()
```



We got good results with CatBosst. We can apply this algorithm on the whole train data set to fit the algorithm. Then make predictions on the test set.

## Predictions on Test Set with CatBoost Algorithm

In this section I applied CatBoost algorithm with parameters of 2000 iteration and 5 folds for cross validation. I made predictions on the test set and saved this probabilities for each client into a csv file. Submitting this csv file (submission\_yagmur\_rigo2.csv) to Kaggle, my score was 0.75399. The screenshot is provided in the attachments (Rapor.doc).



In [ ]:

```

1 import pandas as pd
2 import catboost as cat
3 from sklearn.preprocessing import LabelEncoder
4 import os
5 import pickle
6 from sklearn.metrics import roc_auc_score
7
8
9 directory = 'C:\\CreditRiskProject\\Datasets'
10 train_set = pd.read_csv(directory+ '\\train_set_final.csv', index_col=0)
11 test_set = pd.read_csv(directory+ '\\test_set_final.csv', index_col=0)
12
13
14 #Remove Unnecessary colums
15 def removenonuniquecol(dataset):
16     dropcols = [col for col in dataset.columns if dataset[col].nunique(dropna=True)==1]
17     print ('Removing columns: ',dropcols)
18     dataset.drop(dropcols,axis= 1,inplace= True,errors= 'ignore')
19     return dropcols
20
21 #For Label encoding
22 def labelencoder(dataset):
23     objectlist = dataset.select_dtypes(include=['object']).copy()
24     cat_col = [col for col in dataset.columns if col in objectlist]
25     for col in cat_col:
26         print("Encoding ",col)
27         lbl = LabelEncoder()
28         dataset[col].fillna(-999)
29         lbl.fit(list(dataset[col].values.astype('str')))
30         dataset[col] = lbl.transform(list(dataset[col].values.astype('str')))
31     return cat_col
32
33 #####Trainset-Testset preperation#####
34
35 def TrainPrep(Datasetname):
36     removedCols = removenonuniquecol(Datasetname)
37     Predictors = [col for col in Datasetname]
38     Predictors = [col for col in Predictors if col not in removedCols]
39     Predictors = [col for col in Predictors if col not in ['TARGET']]
40     DfLabel = Datasetname['TARGET']
41     encodedList = labelencoder(Datasetname)
42     return Predictors,Datasetname, DfLabel, removedCols, encodedList
43
44
45 def TestPrep(Datasetname):
46     encodedList = labelencoder(Datasetname)
47     return Datasetname, encodedList
48
49
50
51 Predictors,Df_train, label_train, removedCols, encodedList = TrainPrep(train_set)
52 Df_test, encodedList_test = TestPrep(test_set)
53
54 #####Training Classifier#####
55
56
57 def catboosttrainer(X,y,features,initparam,modelname,modelpath,docpath,cvfold = 5):
58     print ("searching for optimal iteration count...")
59     trainpool = cat.Pool(X[features],y)

```

```

60 cvresult = cat.cv(params= initparam, fold_count=cvfold, pool=trainpool,stratified
61 initparam['iterations'] = (len(cvresult)) - (initparam['od_wait']+1)
62 del initparam['od_wait']
63 del initparam['od_type']
64 print ("optimal iteration count is ", initparam['iterations'])
65 print ("fitting model...")
66 clf = cat.CatBoostClassifier(** initparam)
67 clf.fit(trainpool)
68 imp = clf.get_feature_importance(trainpool,fstr_type='FeatureImportance')
69 dfimp = pd.DataFrame(imp,columns = ['CatBoostImportance'])
70 dfimp.insert(0,column='Feature', value=features)
71 dfimp = dfimp.sort_values(['CatBoostImportance','Feature'], ascending= False)
72 xlsxpath = os.path.join(docpath,modelname+".xlsx")
73 dfimp.to_excel(xlsxpath)
74 print ("pickling model...")
75 picklepath = os.path.join(modelpath,modelname)
76 with open(picklepath,'wb') as fout:
77     pickle.dump(clf, fout)
78     return cvresult,clf,initparam,dfimp
79
80
81
82 modelpath = 'C:\\CreditRiskProject'
83 docpath = 'C:\\CreditRiskProject'
84 CatBoostParam = { 'iterations': 2000, 'od_type': 'Iter', 'od_wait': 100,'loss_function'
85
86 cvresult,clf,initparam,dfimp = catboosttrainer(Df_train,label_train,Predictors,CatBoos
87
88 predictions = clf.predict_proba(Df_test)[:,:1]
89
90
91 #READING SAMPLE SUBMISSION FILE
92
93 sample = pd.read_csv(directory+'\\sample_submission.csv')
94 sample['TARGET']=predictions
95 #CREATING SUBMISSION FILE
96 sample.to_csv(directory+ '\\submission_yagmur_rigo2.csv',index=False)

```

## Predictions on Test Set with Random Forest Algorithm

In this section I applied Random Forest algorithm with parameters of 100 iterations. I made predictions on the test set and saved this probabilities for each client into a csv file. Submitting this csv file (submission\_yagmur\_rigo2\_RandomForest.csv) to Kaggle, my score was 0.66598. The screenshot is provided in the attachments (Rapor.doc).

In [ ]:

```

1 import pandas as pd
2 from sklearn.preprocessing import LabelEncoder
3 from sklearn.ensemble import RandomForestRegressor
4
5 directory = 'C:\\\\CreditRiskProject\\\\Datasets'
6 train_set = pd.read_csv(directory+ '\\\\train_set_final.csv', index_col=0)
7 test_set = pd.read_csv(directory+ '\\\\test_set_final.csv', index_col=0)
8
9 #Remove Unnecessary columns
10 def removenonuniquecol(dataset):
11     dropcols = [col for col in dataset.columns if dataset[col].nunique(dropna=True)==1]
12     print ('Removing columns: ',dropcols)
13     dataset.drop(dropcols,axis= 1,inplace= True,errors= 'ignore')
14     return dropcols
15
16 #For Label encoding
17 def labelencoder(dataset):
18     objectlist = dataset.select_dtypes(include=['object']).copy()
19     cat_col = [col for col in dataset.columns if col in objectlist]
20     for col in cat_col:
21         print("Encoding ",col)
22         lbl = LabelEncoder()
23         dataset[col].fillna(-999)
24         lbl.fit(list(dataset[col].values.astype('str')))
25         dataset[col] = lbl.transform(list(dataset[col].values.astype('str')))
26     return cat_col
27
28 #####Trainset-Testset preperation#####
29
30 def TrainPrep(Datasetname):
31     removedCols = removenonuniquecol(Datasetname)
32     Predictors = [col for col in Datasetname]
33     Predictors = [col for col in Predictors if col not in removedCols]
34     Predictors = [col for col in Predictors if col not in ['TARGET']]
35     DfLabel = Datasetname['TARGET']
36     encodedList = labelencoder(Datasetname)
37     return Predictors,Datasetname, DfLabel, removedCols, encodedList
38
39
40 def TestPrep(Datasetname):
41     encodedList = labelencoder(Datasetname)
42     return Datasetname, encodedList
43
44 #Prepare train and test sets for the model
45 Predictors,Df_train, label_train, removedCols, encodedList = TrainPrep(train_set)
46 Df_test, encodedList_test = TestPrep(test_set)
47
48 X_train = Df_train[Predictors]
49 Y_tarin = label_train
50 x_test = Df_test
51
52 #fill na values with the mean of the column
53 X_train = X_train.fillna(X_train.mean())
54
55 #Fit the RF model:
56 regressor = RandomForestRegressor(n_estimators = 100, random_state = 0)
57 regressor.fit(X_train,Y_tarin)
58
59 #Predict y_test

```

```
60 x_test = x_test.fillna(x_test.mean())
61 y_test = regressor.predict(x_test)
62
63 #READING SAMPLE SUBMISSION FILE
64 sample = pd.read_csv(directory+'\\sample_submission.csv')
65 sample['TARGET']=y_test
66 #CREATING SUBMISSION FILE
67 sample.to_csv(directory+ '\\submission_yagmur_rigo2_RandomForest.csv',index=False)
68
```

In [ ]:

1