

Home Credit Default Risk

1. Business/Real-world Problem

1.1 Problem Statement

Build a model which can predict how capable each applicant is of repaying a loan.

1.2 Source/Useful Links

<https://www.kaggle.com/c/home-credit-default-risk/team> (<https://www.kaggle.com/c/home-credit-default-risk/team>)

1.3 Real-world/Business objectives and constraints.

- No strict latency constraints.
- Interpretability is important.
- Mis-Classification cost is very high.

2. Machine Learning Problem

2.1 Data Description

application_{train|test}.csv

- This is the main table, broken into two files for Train (with TARGET) and Test (without TARGET).
- Static data for all applications. One row represents one loan in our data sample.

bureau.csv

- All client's previous credits provided by other financial institutions that were reported to Credit Bureau (for clients who have a loan in our sample).
- For every loan in our sample, there are as many rows as number of credits the client had in Credit Bureau before the application date.

bureau_balance.csv

- Monthly balances of previous credits in Credit Bureau.
- This table has one row for each month of history of every previous credit reported to Credit Bureau – i.e the table has (#loans in sample * # of relative previous credits * # of months where we have some history observable for the previous credits) rows.

POS_CASH_balance.csv

- Monthly balance snapshots of previous POS (point of sales) and cash loans that the applicant had with Home Credit.
- This table has one row for each month of history of every previous credit in Home Credit (consumer credit and cash loans) related to loans in our sample – i.e. the table has (#loans in sample * # of relative previous credits * # of months in which we have some history observable for the previous credits) rows.

credit_card_balance.csv

- Monthly balance snapshots of previous credit cards that the applicant has with Home Credit.
- This table has one row for each month of history of every previous credit in Home Credit (consumer credit and cash loans) related to loans in our sample – i.e. the table has (#loans in sample * # of relative previous credit cards * # of months where we have some history observable for the previous credit card) rows.

previous_application.csv

- All previous applications for Home Credit loans of clients who have loans in our sample.
- There is one row for each previous application related to loans in our data sample.

installments_payments.csv

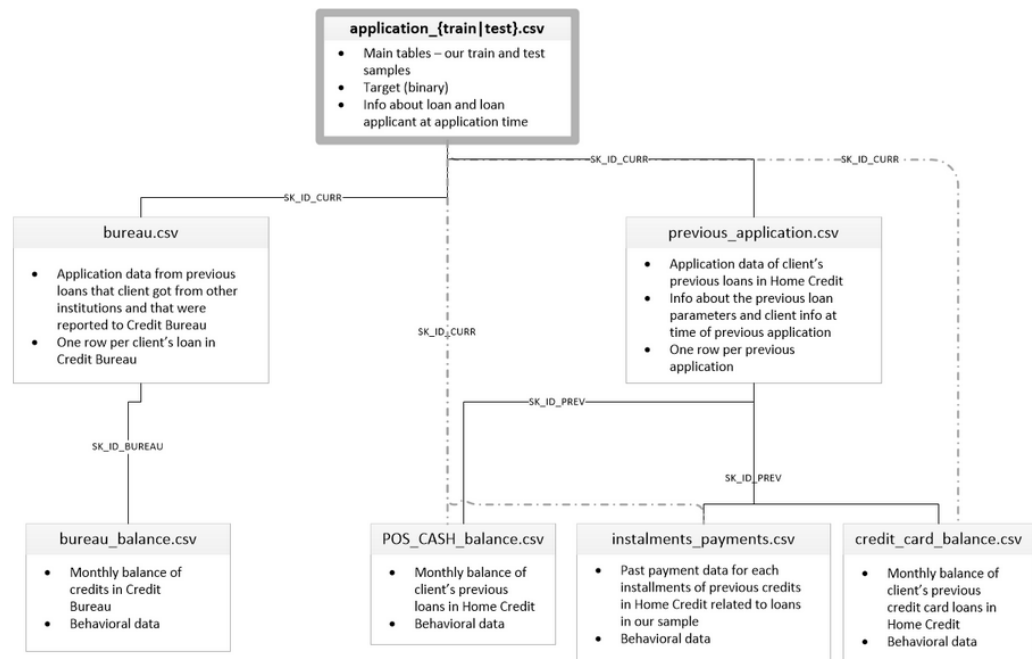
- Repayment history for the previously disturbed credits in Home Credit related to the loans in our sample.
- There is a) one row for every payment that was made plus b) one row each for missed payment.
- One row is equivalent to one payment of one installment OR one installment corresponding to one payment of one previous Home Credit credit related to loans in our sample.

HomeCredit_columns_description.csv

- This file contains descriptions for the columns in the various data files.

```
In [0]: 1 from IPython.display import Image
        2 Image(filename='flow.png')
```

Out[1]:



2.2 Type of Machine Learning Problem

It is a binary class classification problem where 0 means person is capable of repaying a loan and 1 means person is not capable of repaying a loan.

2.3 Performance Metric

- AUC
- Confusion matrix

3. Download the dataset

```
In [1]: 1 !wget --header="Host: storage.googleapis.com" --header="User-Agent: Mozilla/5.0 (X11; Linux x86_64; rv:68.0) Gecko/20100101 Firefox/68.0" https://storage.googleapis.com/kaggle-competitions-data/kaggle-v2/9120/860599/bundle/archive.zip?GoogleAccessId=web-data@kaggle-161607.iam.gserviceaccount.com&Expires=1580916655&Signature=N8bGIC9z3FSFcDpTb9x%2FET5xG5B%2F0dlkxUKQ6hWVwr%2FVyoXki6x%2BBtVlyvckcPiP5nLDP5KvqrmLxBkKcIVkYK8JQTe%2BGY6fL2YMQlVht5Fldlw06cnJikB8Yp2vyuC2RZYt5QSXva6MID%2BpvWI8njIYHMMTpffTXVE8W5ubeaRlt3fUYb2ar8nHUKAfK4fwCjacF14HvGcDYCwKskN3o9YWI8mY2%2BG1UiQkVDFWsST4LM6hsEKd87xHo09000K0WDjFIwGCfpw%2BYgbG0SDn8uZb%2FF4vMKdBZK4vMbpLIEl05Ds1T2%2F0d2I4mkICjVa8ZqFpYIv6jA0tLmdAkFQ%3D%3D&response-content-disposition=attachment%3B+filename%3Dhome-credit-default-risk.zip
--2020-02-03 05:19:17-- https://storage.googleapis.com/kaggle-competitions-data/kaggle-v2/9120/860599/bundle/archive.zip?GoogleAccessId=web-data@kaggle-161607.iam.gserviceaccount.com&Expires=1580916655&Signature=N8bGIC9z3FSFcDpTb9x%2FET5xG5B%2F0dlkxUKQ6hWVwr%2FVyoXki6x%2BBtVlyvckcPiP5nLDP5KvqrmLxBkKcIVkYK8JQTe%2BGY6fL2YMQlVht5Fldlw06cnJikB8Yp2vyuC2RZYt5QSXva6MID%2BpvWI8njIYHMMTpffTXVE8W5ubeaRlt3fUYb2ar8nHUKAfK4fwCjacF14HvGcDYCwKskN3o9YWI8mY2%2BG1UiQkVDFWsST4LM6hsEKd87xHo09000K0WDjFIwGCfpw%2BYgbG0SDn8uZb%2FF4vMKdBZK4vMbpLIEl05Ds1T2%2F0d2I4mkICjVa8ZqFpYIv6jA0tLmdAkFQ%3D%3D&response-content-disposition=attachment%3B+filename%3Dhome-credit-default-risk.zip
Resolving storage.googleapis.com (storage.googleapis.com)... 108.177.97.128, 2404:6800:4008:c04::80
Connecting to storage.googleapis.com (storage.googleapis.com)|108.177.97.128|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 721616255 (688M) [application/zip]
Saving to: 'home-credit-default-risk.zip'

home-credit-default 100%[=====] 688.19M 45.4MB/s in 14s

2020-02-03 05:19:32 (47.5 MB/s) - 'home-credit-default-risk.zip' saved [721616255/721616255]
```

```
In [2]: 1 !unzip /content/home-credit-default-risk.zip
```

```
Archive: /content/home-credit-default-risk.zip
  inflating: HomeCredit_columns_description.csv
  inflating: POS_CASH_balance.csv
  inflating: application_test.csv
  inflating: application_train.csv
  inflating: bureau.csv
  inflating: bureau_balance.csv
  inflating: credit_card_balance.csv
  inflating: installments_payments.csv
  inflating: previous_application.csv
  inflating: sample_submission.csv
```

4. Include packages

```
In [0]: 1 import warnings
2 warnings.filterwarnings("ignore")
3 import shutil
4 import os
5 import pandas as pd
6 import matplotlib
7 matplotlib.use('nbAgg')
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import numpy as np
11 import pickle
12 from sklearn.manifold import TSNE
13 from sklearn import preprocessing
14 import pandas as pd
15 from multiprocessing import Process# this is used for multithreading
16 import multiprocessing
17 import codecs# this is used for file operations
18 import random as r
19 from xgboost import XGBClassifier
20 from sklearn.model_selection import RandomizedSearchCV
21 from sklearn.tree import DecisionTreeClassifier
22 from sklearn.calibration import CalibratedClassifierCV
23 from sklearn.neighbors import KNeighborsClassifier
24 from sklearn.metrics import log_loss
25 from sklearn.metrics import confusion_matrix
26 from sklearn.model_selection import train_test_split
27 from sklearn.linear_model import LogisticRegression
28 from sklearn.ensemble import RandomForestClassifier
29 import math
30 from math import log
31 %matplotlib inline
32 import gc
33 import pickle
34 from lightgbm import LGBMClassifier
35 from sklearn.metrics import roc_auc_score, roc_curve, auc
36 from sklearn.model_selection import KFold, StratifiedKFold
37 import seaborn as sns
38 import warnings
39 warnings.simplefilter(action='ignore', category=FutureWarning)
40
```

5. Perform EDA

```
In [0]: 1 Df_application_train=pd.read_csv("application_train.csv")
```

In [0]: 1 Df application train

Out[6]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_I
0	100002	1	Cash loans	M	N	
1	100003	0	Cash loans	F	N	
2	100004	0	Revolving loans	M	Y	
3	100006	0	Cash loans	F	N	
4	100007	0	Cash loans	M	N	
...
307506	456251	0	Cash loans	M	N	
307507	456252	0	Cash loans	F	N	
307508	456253	0	Cash loans	F	N	
307509	456254	1	Cash loans	F	N	
307510	456255	0	Cash loans	F	N	

307511 rows × 122 columns

In [0]: 1 Df application train.TARGET.value counts()

Out[7]: 0 282686
1 24825
Name: TARGET, dtype: int64

The dataset is highly imbalanced. There are more number of person in the dataset who are capable of giving the loan.

In [0]: 1 Df application train.columns

Out[133]: Index(['SK_ID_CURR', 'TARGET', 'NAME_CONTRACT_TYPE', 'CODE_GENDER',
'FLAG_OWN_CAR', 'FLAG_OWN_REALTY', 'CNT_CHILDREN', 'AMT_INCOME_TOTA
L',
'AMT_CREDIT', 'AMT_ANNUITY',
'...',
'FLAG_DOCUMENT_18', 'FLAG_DOCUMENT_19', 'FLAG_DOCUMENT_20',
'FLAG_DOCUMENT_21', 'AMT_REQ_CREDIT_BUREAU_HOUR',
'AMT_REQ_CREDIT_BUREAU_DAY', 'AMT_REQ_CREDIT_BUREAU_WEEK',
'AMT_REQ_CREDIT_BUREAU_MON', 'AMT_REQ_CREDIT_BUREAU_QRT',
'AMT_REQ_CREDIT_BUREAU_YEAR'],
dtype='object', length=122)

In [0]: 1 Df application train.NAME CONTRACT TYPE.unique()

Out[134]: array(['Cash loans', 'Revolving loans'], dtype=object)

In [0]: 1 Df application test=nd.read_csv("application test.csv")

In [0]: 1 Df_application_test.columns

Out[11]: Index(['SK_ID_CURR', 'NAME_CONTRACT_TYPE', 'CODE_GENDER', 'FLAG_OWN_CAR',
'FLAG_OWN_REALTY', 'CNT_CHILDREN', 'AMT_INCOME_TOTAL', 'AMT_CREDIT',
,
'AMT_ANNUITY', 'AMT_GOODS_PRICE',
,
'FLAG_DOCUMENT_18', 'FLAG_DOCUMENT_19', 'FLAG_DOCUMENT_20',
'FLAG_DOCUMENT_21', 'AMT_REQ_CREDIT_BUREAU_HOUR',
'AMT_REQ_CREDIT_BUREAU_DAY', 'AMT_REQ_CREDIT_BUREAU_WEEK',
'AMT_REQ_CREDIT_BUREAU_MON', 'AMT_REQ_CREDIT_BUREAU_QRT',
'AMT_REQ_CREDIT_BUREAU_YEAR'],
dtype='object', length=121)

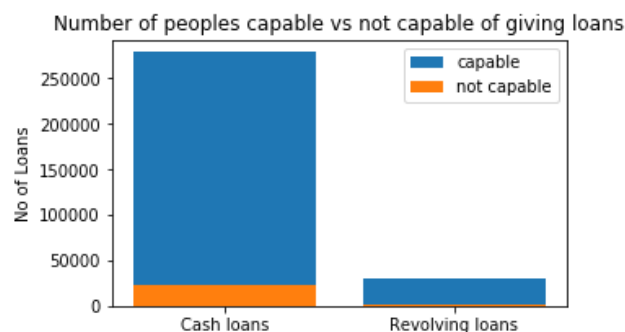
Function to plot bar graph

```
In [0]: 1 def stack_plot(data, xtick, col2, col3='total'):
2     ind = np.arange(data.shape[0])
3     %matplotlib inline
4     plt.figure(figsize=(10,3))
5     p1 = plt.bar(ind, data[col3].values)
6     p2 = plt.bar(ind, data[col2].values)
7     plt.ylabel('No of Loans')
8     plt.title('Number of peoples capable vs not capable of giving loans')
9     plt.xticks(ind, list(data[xtick].values), rotation=-90)
10    plt.legend((p1[0], p2[0]), ('capable', 'not capable'))
11    plt.show()
12    def univariate_barplots(data, col1, col2, top=False):
13        # Count number of zeros in dataframe python: https://stackoverflow.
14        temp = pd.DataFrame(Df_application_train.groupby(col1)[col2].agg('la
15
16        # Pandas dataframe grouby count: https://stackoverflow.com/a/193855
17        temp['total'] = pd.DataFrame(Df_application_train.groupby(col1)[col
18
19        temp['Avg'] = pd.DataFrame(Df_application_train.groupby(col1)[col2]
20
21        temp.sort_values(by=['total'], inplace=True, ascending=False)
22
23        if top:
24            temp = temp[0:top]
25
26        stack_plot(temp, xtick=col1, col2=col2, col3='total')
27
```

6. Analysis on application.csv file

6.1 Identification if loan is cash or revolving

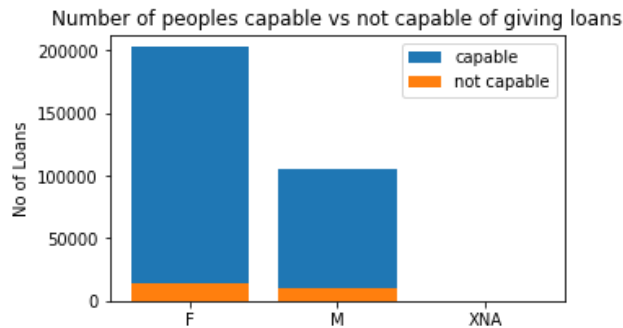
In [0]: 1 univariate_barplots(Df_application_train, 'NAME_CONTRACT_TYPE', 'TARGET')



observation: Many people took cash loans.

6.2 Gender of the client

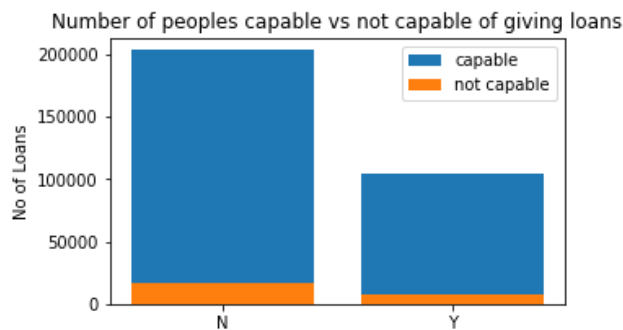
```
In [0]: 1 univariate_barplots(Df_application_train, 'CODE_GENDER', 'TARGET', to
```



observation : Female took more no of loans then males.

6.3 Flag if the client owns a car

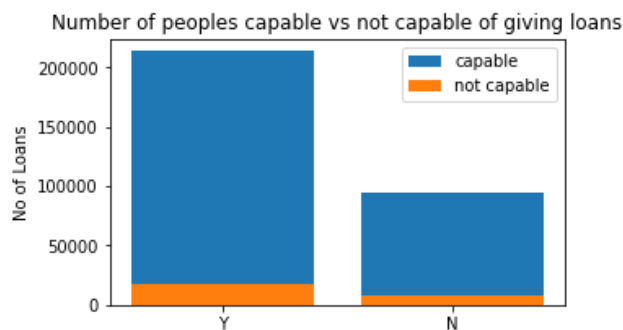
```
In [0]: 1 univariate_barplots(Df_application_train, 'FLAG_OWN_CAR', 'TARGET', to
```



observation : People who own a car took less no of loan.

6.4 Flag if client owns a house or flat

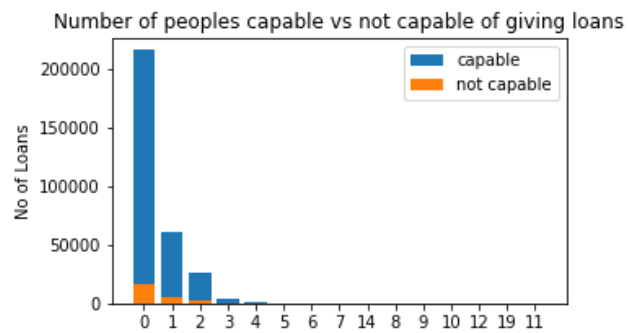
```
In [0]: 1 univariate_barplots(Df_application_train, 'FLAG_OWN_REALTY', 'TARGET', to
```



observation : People who owns a house or flat takes more no of loans.

6.5 Number of children the client has

```
In [0]: 1 univariate_barplots(Df_application_train, 'CNT CHILDREN', 'TARGET', tc
```



observation : Clients having no child take more no of loans and are also more defaulters.

Function to plot pie chart

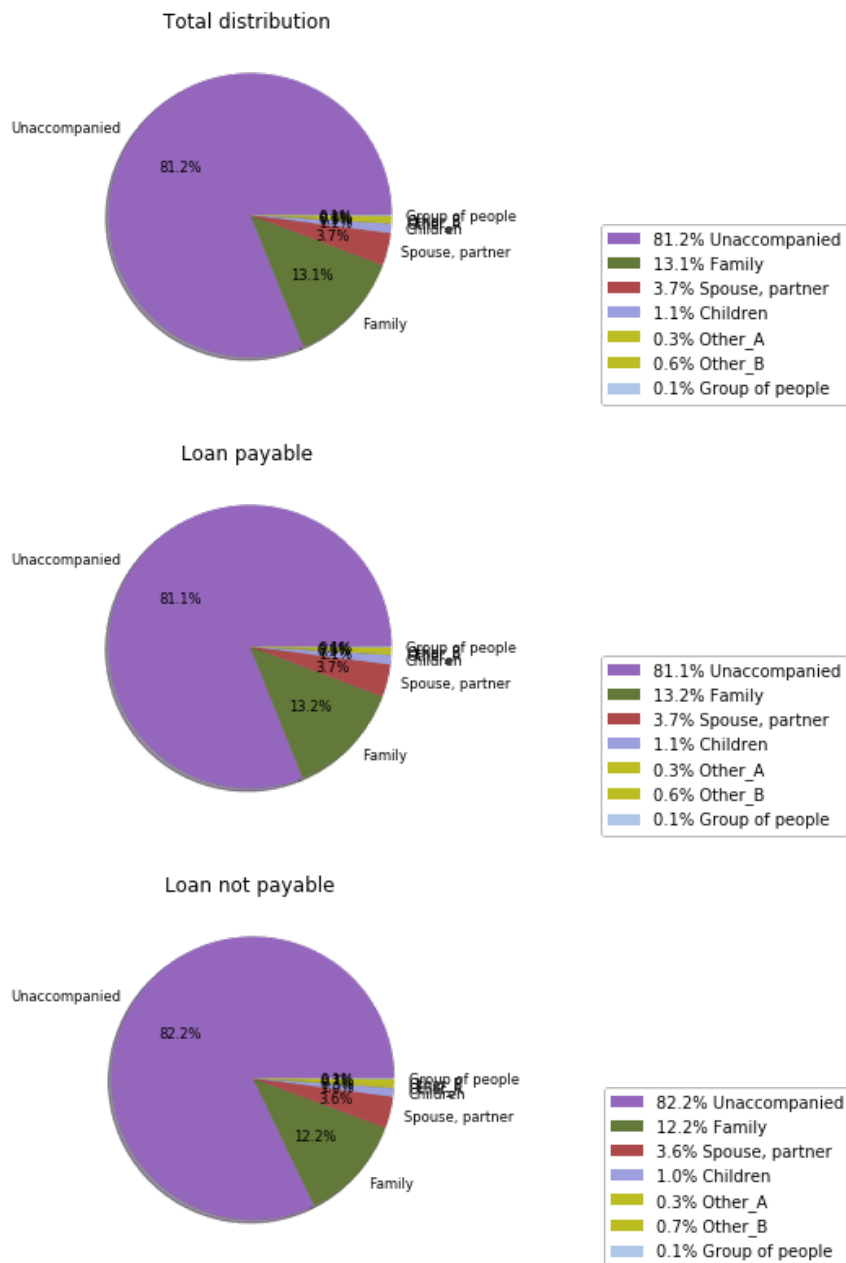
```

In [0]: 1 def labelupdate(autotexts,labels):
2         l=[]
3         for i in range(len(labels)):
4             l.append(str(autotexts[i]._text) + ' ' + str(labels[i]))
5         return(l)
6
7 def plotpiechart(s):
8     %matplotlib inline
9     import random
10    labels=list(Df_application_train[s].unique())
11    labels=[x for x in labels if (str(x)!='nan')]
12    fracs1=[]
13    fracs2=[]
14    fracs3=[]
15    for x in labels:
16        count1 = sum(Df_application_train[s]==x)
17        fracs1.append(count1)
18        count2 = sum(Df_application_train[Df_application_train['TARGET']==0]==x)
19        fracs2.append(count2)
20        count3 = sum(Df_application_train[Df_application_train['TARGET']==1]==x)
21        fracs3.append(count3)
22    colors=["#9c9ede", "#7375b5", "#4a5584", "#cedb9c", "#b5cf6b", "#8ca2c2",
23            "#e7cb94", "#e7ba52", "#bd9e39", "#8c6d31", "#e7969c", "#d6616b",
24            "#843c39", "#de9ed6", "#ce6dbd", "#a55194", "#7b4173", "#1f77b4",
25            "#ff7f0e", "#ffbb78", "#2ca02c", "#98df8a", "#d62728", "#ff9896",
26            "#c5b0d5", "#8c564b", "#c49c94", "#e377c2", "#f7b6d2", "#7f7f7f",
27            "#bcbd22", "#dbdb8d", "#17becf", "#9edae5", "#1f77b4", "#ff7f0e",
28            "#d62728", "#9467bd", "#8c564b", "#e377c2", "#7f7f7f", "#bcbd22"]
29    colors=random.sample(colors, len(labels))
30    plt.title("Total distribution ")
31
32    patches,texts,autotexts=plt.pie(fracs1, labels=labels, colors=colors,
33    l=labelupdate(autotexts,labels)
34    plt.legend(patches,l, bbox_to_anchor=(1,0.5), loc="best", fontsize=16)
35    plt.show()
36    plt.title("Loan payable ")
37    patches,texts,autotexts=plt.pie(fracs2, labels=labels, colors=colors,
38    l=labelupdate(autotexts,labels)
39    plt.legend(patches,l, bbox_to_anchor=(1,0.5), loc="best", fontsize=16)
40    plt.show()
41    plt.title("Loan not payable ")
42    patches,texts,autotexts=plt.pie(fracs3, labels=labels, colors=colors,
43    l=labelupdate(autotexts,labels)
44    plt.legend(patches,l, bbox_to_anchor=(1,0.5), loc="best", fontsize=16)
45    plt.show()
46
47

```

6.6 Who was accompanying client when he was applying for the loan

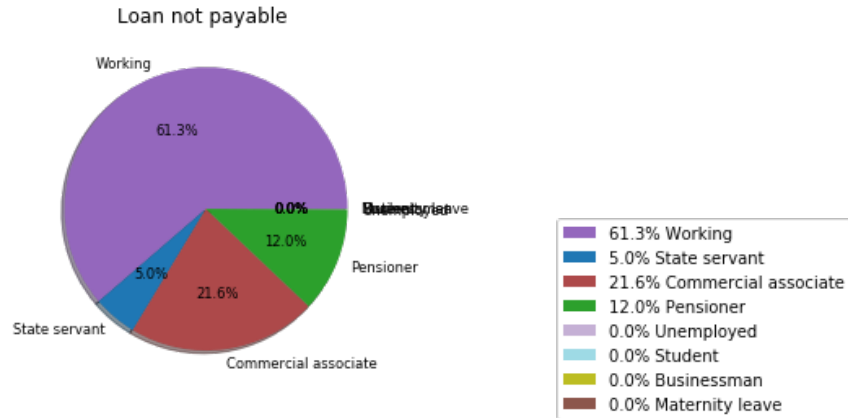
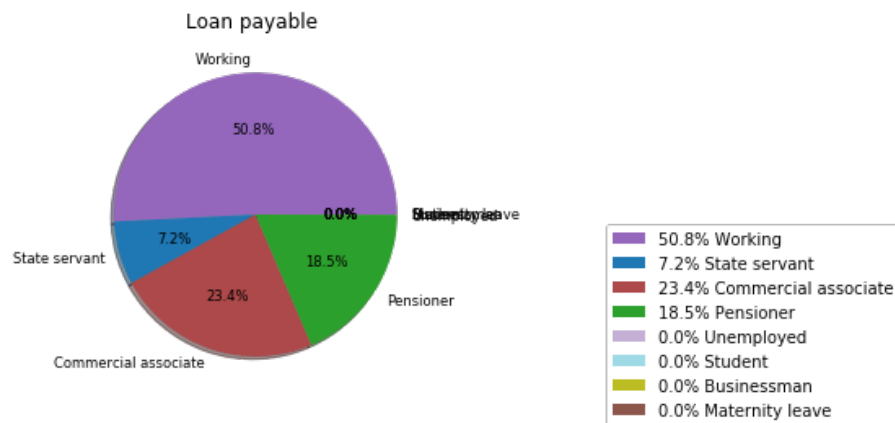
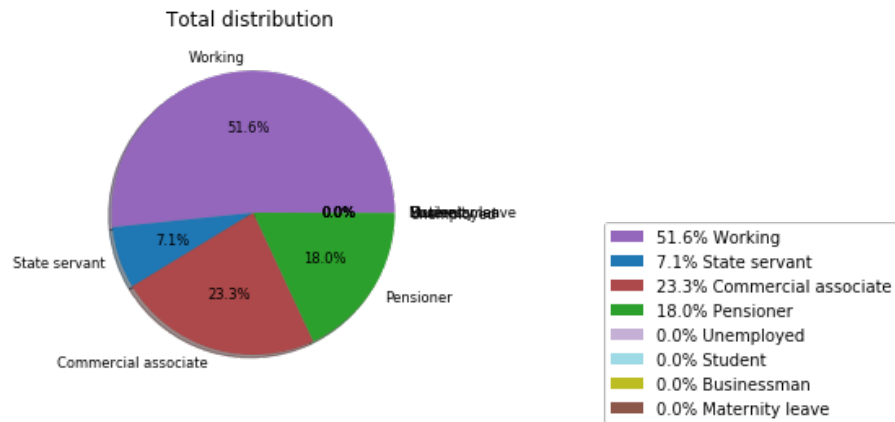
```
In [0]: 1 plotpiechart('NAME TYPE SUITE')
```



observation : Client who came alone took more loan but also were more defaulters.

6.7 Clients income type (Working,State servant,Commercial associate,Pensioner,Unemployed,Student,Businessman,Maternity leave)

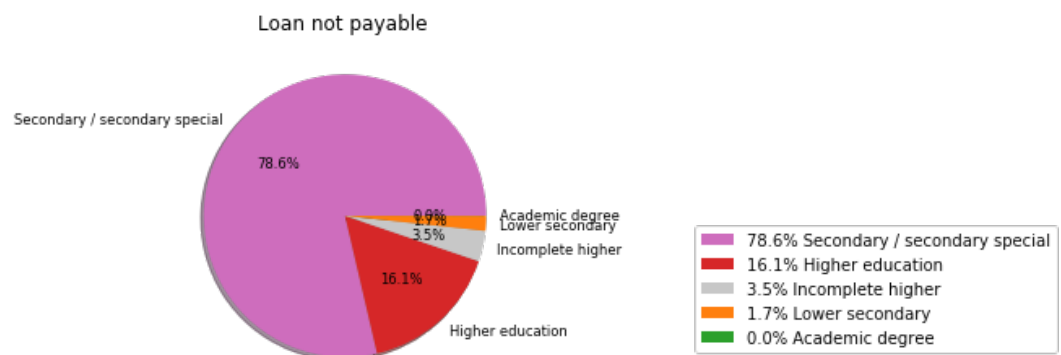
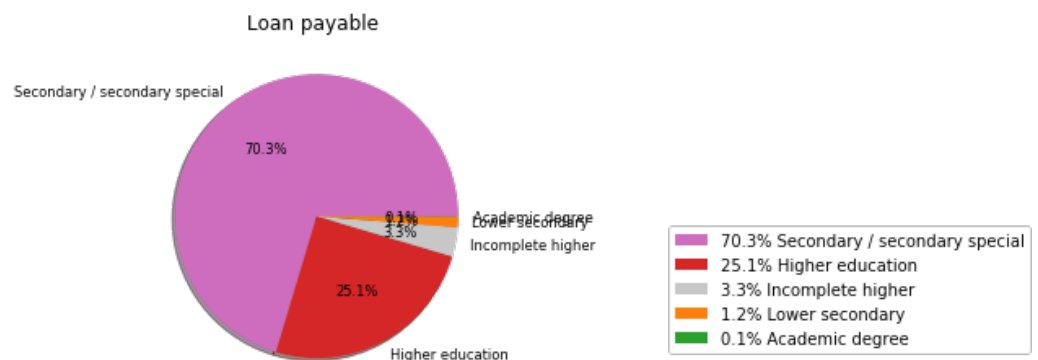
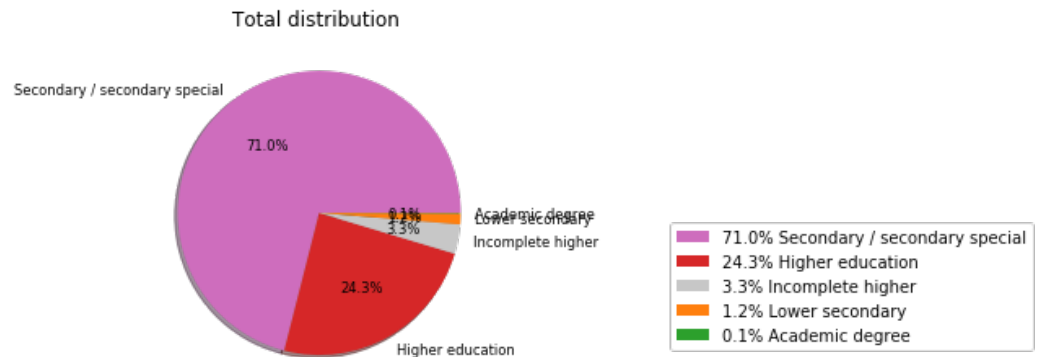
```
In [0]: 1 plotpiechart('NAME INCOME TYPE')
```



observation : clients with income type working take more no of loans.

6.8 Level of highest education the client achieved.

```
In [0]: 1 plotpiechart('NAME EDUCATION TYPE')
```



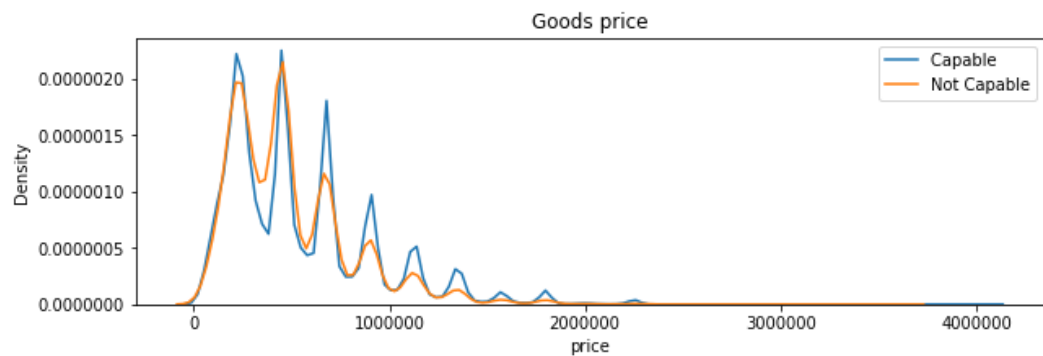
observation : people with highest degree as secondary and secondary special are more no of defaulters compare to people with other degree.

6.9 For consumer loans it is the price of the goods for which the loan is given

```

In [0]: 1 from math import log
        2
        3 count1 = Df_application_train[Df_application_train['TARGET']==0]['AMT_C
        4 count1 = count1.values
        5 count1=[x for x in count1 if (str(x)!='nan')]
        6 count2 = Df_application_train[Df_application_train['TARGET']==1]['AMT_C
        7 count2 = count2.values
        8 count2=[x for x in count2 if (str(x)!='nan')]
        9 plt.figure(figsize=(10,3))
       10 sns.distplot(count1, hist=False, label="Capable ")
       11 sns.distplot(count2, hist=False, label="Not Capable")
       12 plt.title('Goods price')
       13 plt.xlabel('price')
       14 plt.ylabel('Density')
       15 plt.legend()
       16 plt.show()

```



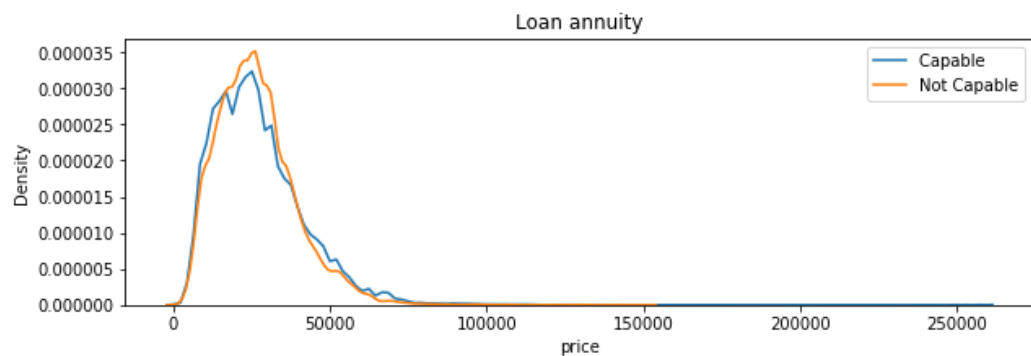
Observation: Most no of loans are given for goods price below 10 lakhs

6.10 Loan annuity

```

In [5]: 1 from math import log
        2
        3 count1 = Df_application_train[Df_application_train['TARGET']==0]['AMT_A
        4 count1 = count1.values
        5 count1=[x for x in count1 if (str(x)!='nan')]
        6 count2 = Df_application_train[Df_application_train['TARGET']==1]['AMT_A
        7 count2 = count2.values
        8 count2=[x for x in count2 if (str(x)!='nan')]
        9 plt.figure(figsize=(10,3))
       10 sns.distplot(count1, hist=False, label="Capable ")
       11 sns.distplot(count2, hist=False, label="Not Capable")
       12 plt.title('Loan annuity')
       13 plt.xlabel('price')
       14 plt.ylabel('Density')
       15 plt.legend()
       16 plt.show()
       17

```



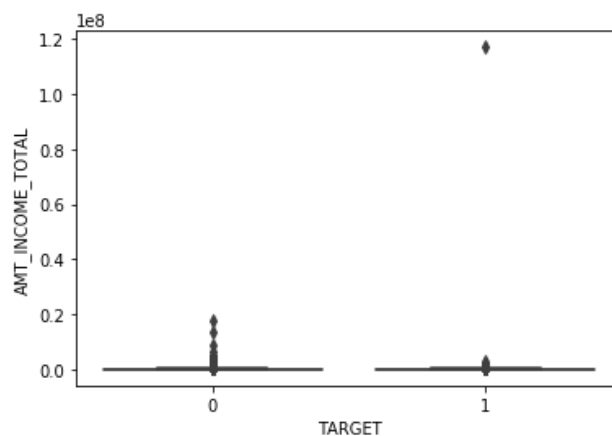
observation : Most people pay annuity below 50000 for the credit loan

6.11 : Income of the client

```

In [0]: 1 sns.boxplot(x='TARGET', y='AMT_INCOME_TOTAL', data=Df_application_train)
        2 plt.show()

```



We can't say anything here let's check the percentile

```
In [0]: 1 v=Df_application_train['AMT_INCOME_TOTAL'].values
        2 v=np.sort(v)
        3 for i in range(0,100,10):
        4     print(i,'percentile value is',v[int(len(v)*float(i)/100)])
        5 print(100,'percentile value is'.v[-1])

0 percentile value is 25650.0
10 percentile value is 81000.0
20 percentile value is 99000.0
30 percentile value is 112500.0
40 percentile value is 135000.0
50 percentile value is 147150.0
60 percentile value is 162000.0
70 percentile value is 180000.0
80 percentile value is 225000.0
90 percentile value is 270000.0
100 percentile value is 117000000.0
```

```
In [0]: 1 v=Df_application_train['AMT_INCOME_TOTAL'].values
        2 v=np.sort(v)
        3 for i in range(90,100,1):
        4     print(i,'percentile value is',v[int(len(v)*float(i)/100)])
        5 print(100,'percentile value is'.v[-1])

90 percentile value is 270000.0
91 percentile value is 270000.0
92 percentile value is 292500.0
93 percentile value is 315000.0
94 percentile value is 315000.0
95 percentile value is 337500.0
96 percentile value is 360000.0
97 percentile value is 382500.0
98 percentile value is 427500.0
99 percentile value is 472500.0
100 percentile value is 117000000.0
```

```
In [0]: 1 v=Df_application_train['AMT_INCOME_TOTAL'].values
        2 v=np.sort(v)
        3 for i in np.arange(0.0,1.0,0.1):
        4     print(99+i,'percentile value is',v[int(len(v)*float(99+i)/100)])
        5 print(100,'percentile value is'.v[-1])

99.0 percentile value is 472500.0
99.1 percentile value is 495000.0
99.2 percentile value is 540000.0
99.3 percentile value is 540000.0
99.4 percentile value is 562500.0
99.5 percentile value is 630000.0
99.6 percentile value is 675000.0
99.7 percentile value is 675000.0
99.8 percentile value is 765000.0
99.9 percentile value is 900000.0
100 percentile value is 117000000.0
```



```
In [0]: 1 v=Df_application_train['AMT_INCOME_TOTAL'].values
2 v=np.sort(v)
3 for i in np.arange(0.90,1.0,0.01):
4     print(99+i,'percentile value is',v[int(len(v)*float(99+i)/100)])
5     print(100.'percentile value is'.v[-1])

99.9 percentile value is 900000.0
99.91 percentile value is 909000.0
99.92 percentile value is 1035000.0
99.93 percentile value is 1125000.0
99.94 percentile value is 1125000.0
99.95 percentile value is 1215000.0
99.96 percentile value is 1350000.0
99.97 percentile value is 1350000.0
99.98 percentile value is 1800000.0
99.99 percentile value is 2250000.0
100 percentile value is 117000000.0
```

```
In [0]: 1 v=Df_application_train['AMT_INCOME_TOTAL'].values
2 v=np.sort(v)
3 for i in np.arange(0.99,0.999,0.001):
4     print(99+i,'percentile value is',v[int(len(v)*float(99+i)/100)])
5     print(100.'percentile value is'.v[-1])

99.99 percentile value is 2250000.0
99.991 percentile value is 2250000.0
99.992 percentile value is 2250000.0
99.993 percentile value is 2250000.0
99.994 percentile value is 2700000.0
99.995 percentile value is 3150000.0
99.996 percentile value is 3600000.0
99.997 percentile value is 3950059.5
99.998 percentile value is 4500000.0
99.999 percentile value is 9000000.0
100 percentile value is 117000000.0
```

```
In [0]: 1 Df_application_train[Df_application_train['AMT_INCOME_TOTAL']>9000000]
```

```
Out[279]:
```

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_I
	12840	114967	1	Cash loans	F	N
	203693	336147	0	Cash loans	M	Y
	246858	385674	0	Cash loans	M	Y

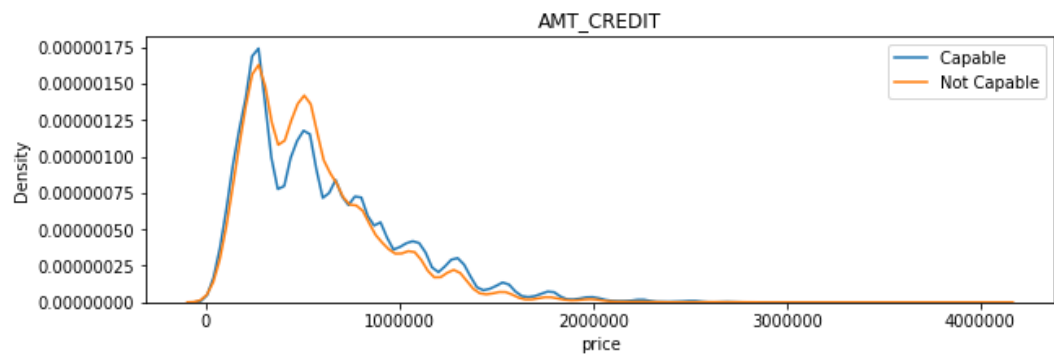
3 rows × 122 columns

6.12 Credit amount of the loan

```

In [0]: 1 from math import log
2 count1 = Df_application_train[Df_application_train['TARGET']==0]['AMT_CREDIT']
3 count1 = count1.values
4 count1=[x for x in count1 if (str(x)!='nan')]
5 count2 = Df_application_train[Df_application_train['TARGET']==1]['AMT_CREDIT']
6 count2 = count2.values
7 count2=[x for x in count2 if (str(x)!='nan')]
8 plt.figure(figsize=(10,3))
9 sns.distplot(count1, hist=False, label="Capable ")
10 sns.distplot(count2, hist=False, label="Not Capable")
11 plt.title('AMT_CREDIT')
12 plt.xlabel('price')
13 plt.ylabel('Density')
14 plt.legend()
15 plt.show()
16

```



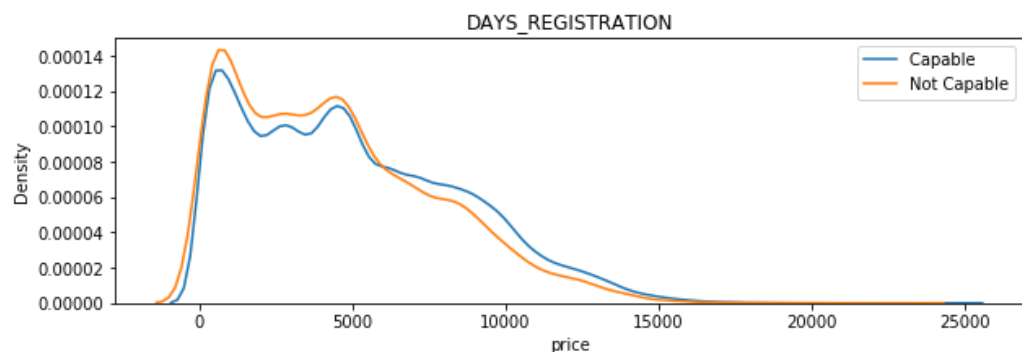
observation : Credit amount of the loan is mostly less then 10 lakhs

6.13 How many days before the application did client change his registration

```

In [0]: 1 from math import log
2 count1 = Df_application_train[Df_application_train['TARGET']==0]['DAYS_REGISTRATION']
3 count1 = count1.values
4 count1=[-x for x in count1 if (str(x)!='nan' and x!=0)]
5 count2 = Df_application_train[Df_application_train['TARGET']==1]['DAYS_REGISTRATION']
6 count2 = count2.values
7 count2=[-x for x in count2 if (str(x)!='nan' and x!=0)]
8 plt.figure(figsize=(10,3))
9 sns.distplot(count1, hist=False, label="Capable ")
10 sns.distplot(count2, hist=False, label="Not Capable")
11 plt.title('DAYS_REGISTRATION')
12 plt.xlabel('price')
13 plt.ylabel('Density')
14 plt.legend()
15 plt.show()
16

```



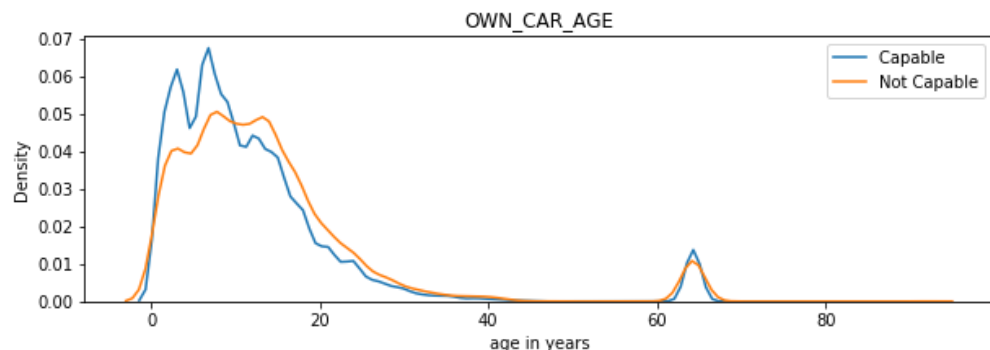
Observation: Most of the applicant changed their registration 15000 days before the application.

6.14 Age of client's car

```
In [0]: 1 Df_application_train['OWN_CAR_AGE'].describe()
```

```
Out[201]: count      104582.000000
mean         12.061091
std          11.944812
min           0.000000
25%           5.000000
50%           9.000000
75%          15.000000
max          91.000000
Name: OWN_CAR_AGE, dtype: float64
```

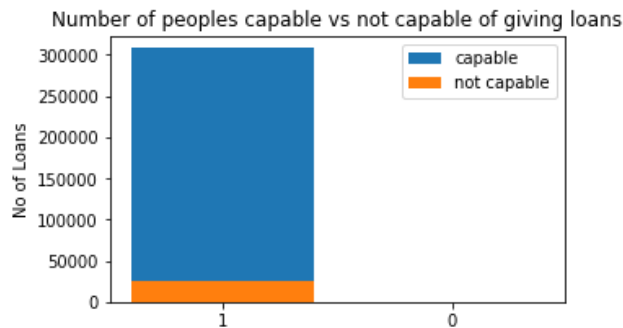
```
In [0]: 1 from math import log
2 count1 = Df_application_train[Df_application_train['TARGET']==0]['OWN_CAR_AGE']
3 count1 = count1.values
4 count1=[x for x in count1 if (str(x)!='nan' and x!=0)]
5 count2 = Df_application_train[Df_application_train['TARGET']==1]['OWN_CAR_AGE']
6 count2 = count2.values
7 count2=[x for x in count2 if (str(x)!='nan' and x!=0)]
8 plt.figure(figsize=(10,3))
9 sns.distplot(count1, hist=False, label="Capable ")
10 sns.distplot(count2, hist=False, label="Not Capable")
11 plt.title('OWN_CAR_AGE')
12 plt.legend()
13 plt.xlabel('age in years')
14 plt.ylabel('Density')
15 plt.show()
```



observation : Most of the cars are in age group between 0-20 and less between 40-60.

6.15 Did client provide mobile phone (1=YES, 0=NO)

```
In [0]: 1 univariate_barplots(Df_application_train, 'FLAG_MOBIL', 'TARGET', top=
```



observation: nearly all applicants have provided their mobile information. Lets check how many applicants are there who have not provided the mobile information.

```
In [0]: 1 c=Df_application_train.FLAG_MOBIL.values
2 co=0
3 for i in c:
4     if(i!=1):
5         co=co+1
6 co
```

Out[245]: 1

```
In [0]: 1 Df_application_train[Df_application_train['FLAG_MOBIL']!=0]
```

Out[8]:

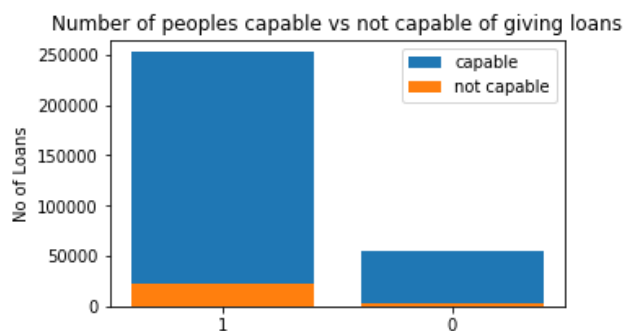
	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_R
15709	118330	0	Cash loans	M	Y	

1 rows × 122 columns

There is only 1 client who did not provide the mobile information.

6.16 Did client provide work phone (1=YES, 0=NO)

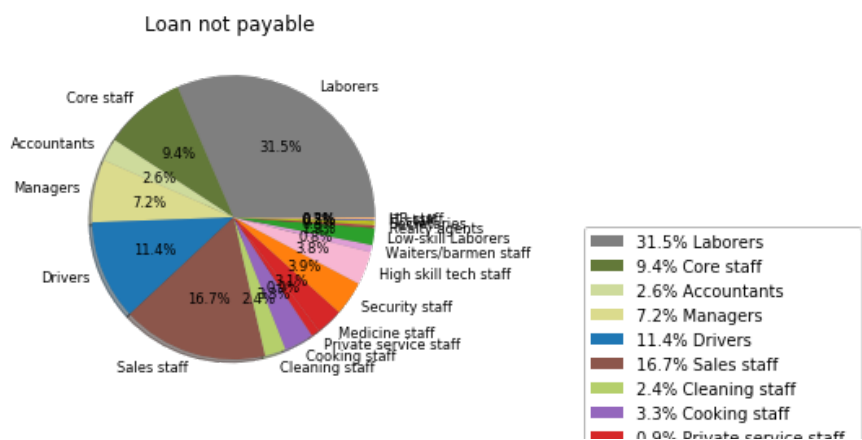
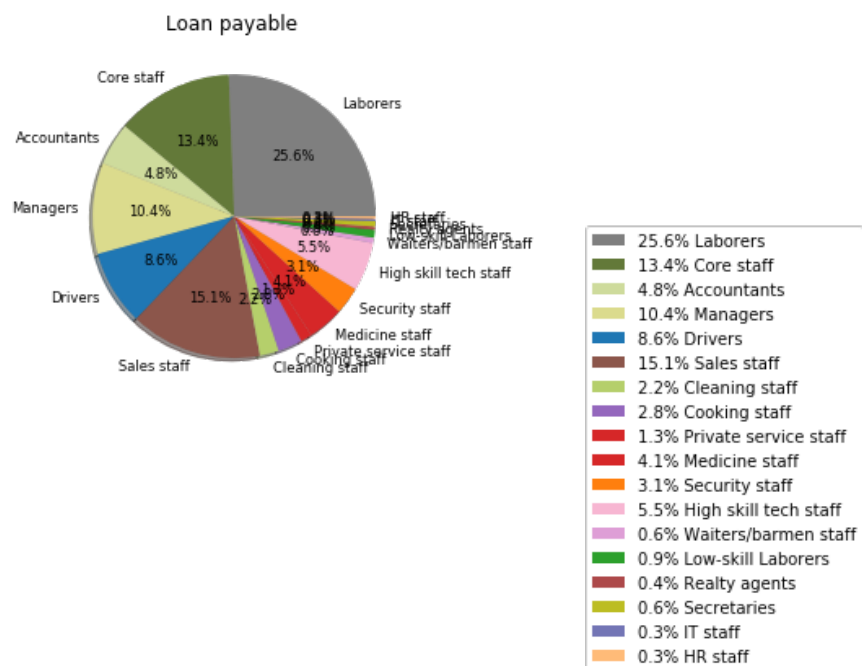
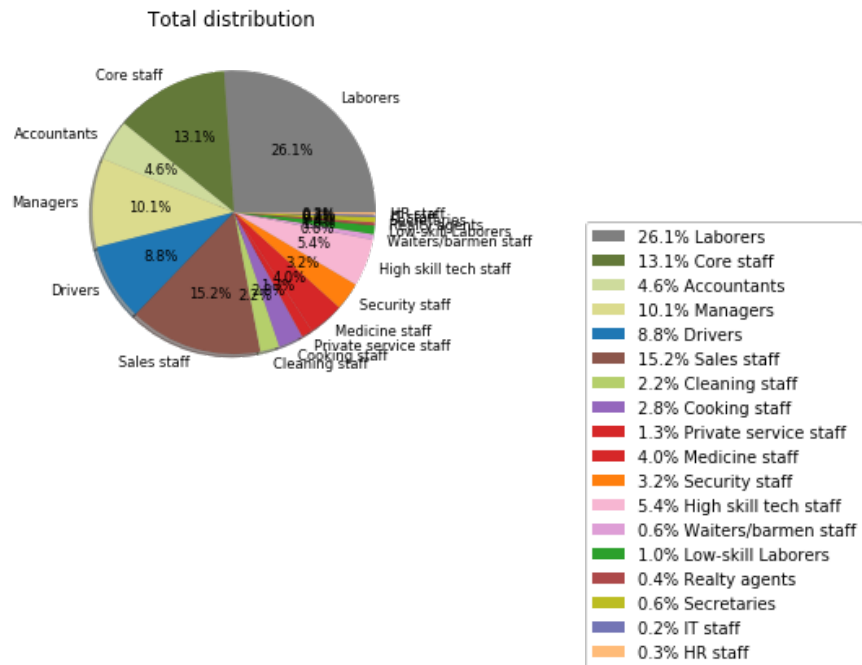
```
In [0]: 1 univariate_barplots(Df_application_train, 'FLAG_EMP_PHONE', 'TARGET',
```



observation : More no of applicants provide work phone.

6.17 What kind of occupation does the client have

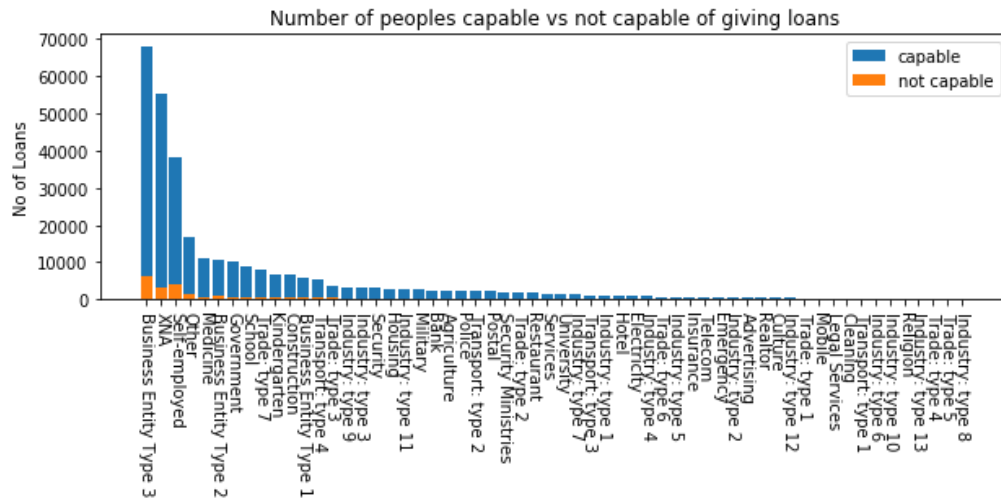
```
In [0]: 1 plotpiechart('OCCUPATION TYPE')
```



observation : Laborers take more no of loans but also are more no of defaulters.

6.18 Type of organization where client works

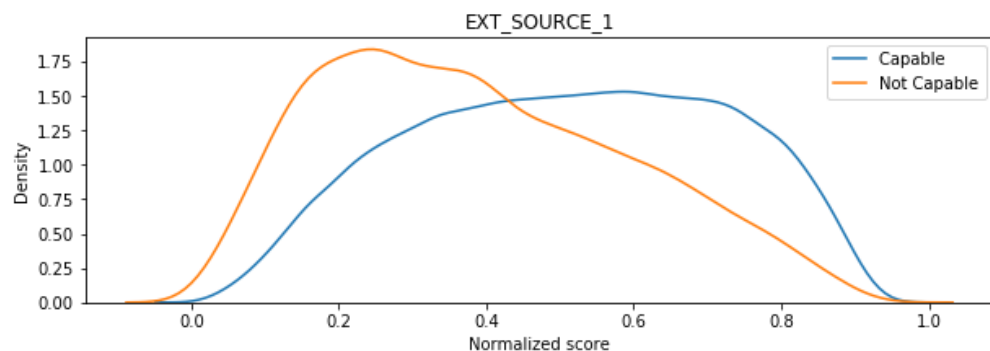
```
In [0]: 1 univariate barplots(Df_application_train, 'ORGANIZATION TYPE', 'TARGET')
```



observation : clients working in business and self employed are more no of defaulters but they also take more no of loan.

6.19 Normalized score from external data source

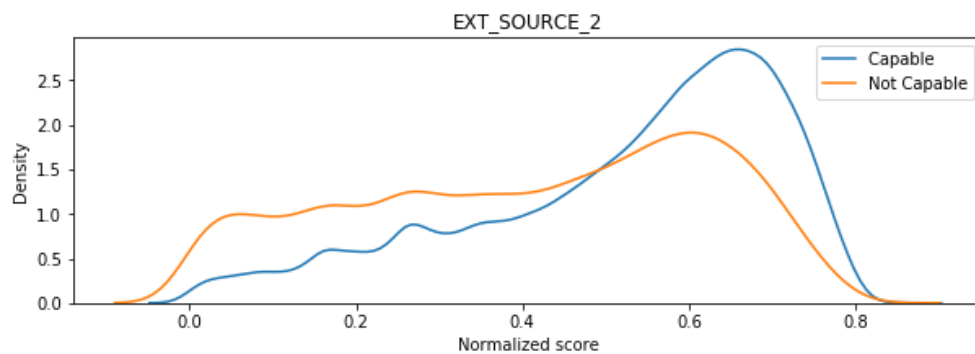
```
In [0]: 1 from math import log
2 count1 = Df_application_train[Df_application_train['TARGET']==0]['EXT_S
3 count1 = count1.values
4 count1=[x for x in count1 if (str(x)!='nan' and x!=0)]
5 count2 = Df_application_train[Df_application_train['TARGET']==1]['EXT_S
6 count2 = count2.values
7 count2=[x for x in count2 if (str(x)!='nan' and x!=0)]
8 plt.figure(figsize=(10,3))
9 sns.distplot(count1, hist=False, label="Capable ")
10 sns.distplot(count2, hist=False, label="Not Capable")
11 plt.title('EXT_SOURCE_1')
12 plt.xlabel('Normalized score')
13 plt.ylabel('Density')
14 plt.legend()
15 plt.show()
```



observation : data is well seperated it will be a important feature.

6.20 Normalized score from external data source

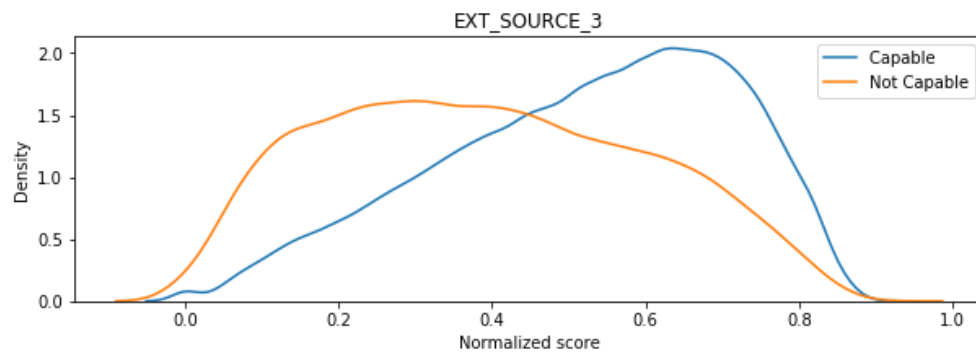
```
In [0]: 1 from math import log
2 count1 = Df_application_train[Df_application_train['TARGET']==0]['EXT_S
3 count1 = count1.values
4 count1=[x for x in count1 if (str(x)!='nan' and x!=0)]
5 count2 = Df_application_train[Df_application_train['TARGET']==1]['EXT_S
6 count2 = count2.values
7 count2=[x for x in count2 if (str(x)!='nan' and x!=0)]
8 plt.figure(figsize=(10,3))
9 sns.distplot(count1, hist=False, label="Capable ")
10 sns.distplot(count2, hist=False, label="Not Capable")
11 plt.title('EXT_SOURCE_2')
12 plt.xlabel('Normalized score')
13 plt.ylabel('Density')
14 plt.legend()
15 plt.show()
```



observation : data is well seperated it will be a important feature

6.21 Normalized score from external data source

```
In [0]: 1 from math import log
2 count1 = Df_application_train[Df_application_train['TARGET']==0]['EXT_S
3 count1 = count1.values
4 count1=[x for x in count1 if (str(x)!='nan' and x!=0)]
5 count2 = Df_application_train[Df_application_train['TARGET']==1]['EXT_S
6 count2 = count2.values
7 count2=[x for x in count2 if (str(x)!='nan' and x!=0)]
8 plt.figure(figsize=(10,3))
9 sns.distplot(count1, hist=False, label="Capable ")
10 sns.distplot(count2, hist=False, label="Not Capable")
11 plt.title('EXT_SOURCE_3')
12 plt.xlabel('Normalized score')
13 plt.ylabel('Density')
14 plt.legend()
15 plt.show()
```



observation : data is well seperated it will be a important feature

Below code will reduce the memory used by the data frame

```
In [18]: 1 np.iinfo(np.int8).max
```

```
Out[18]: 127
```



```

In [0]: 1 # code taken from https://www.kaggle.com/gemartin/load-data-reduce-memory-usage
2 def reduce_mem_usage(df):
3     """ iterate through all the columns of a dataframe and modify the data type
4     to reduce memory usage.
5     """
6     start_mem = df.memory_usage().sum() / 1024**2
7     print('Memory usage of dataframe is {:.2f} MB'.format(start_mem))
8
9     for col in df.columns:
10        col_type = df[col].dtype
11
12        if col_type != object:
13            c_min = df[col].min()
14            c_max = df[col].max()
15            if str(col_type)[:3] == 'int':
16                if c_min > np.iinfo(np.int8).min and c_max < np.iinfo(np.int8).max:
17                    df[col] = df[col].astype(np.int8)
18                elif c_min > np.iinfo(np.int16).min and c_max < np.iinfo(np.int16).max:
19                    df[col] = df[col].astype(np.int16)
20                elif c_min > np.iinfo(np.int32).min and c_max < np.iinfo(np.int32).max:
21                    df[col] = df[col].astype(np.int32)
22                elif c_min > np.iinfo(np.int64).min and c_max < np.iinfo(np.int64).max:
23                    df[col] = df[col].astype(np.int64)
24            else:
25                if c_min > np.finfo(np.float16).min and c_max < np.finfo(np.float16).max:
26                    df[col] = df[col].astype(np.float16)
27                elif c_min > np.finfo(np.float32).min and c_max < np.finfo(np.float32).max:
28                    df[col] = df[col].astype(np.float32)
29                else:
30                    df[col] = df[col].astype(np.float64)
31
32    end_mem = df.memory_usage().sum() / 1024**2
33    print('Memory usage after optimization is: {:.2f} MB'.format(end_mem))
34    print('Decreased by {:.1f}%'.format(100 * (start_mem - end_mem) / start_mem))
35
36    return df
37

```

Feature engineering for application .csv

Remove outliers

```

In [0]: 1 (Df_application_train['DAYS_BIRTH'] / -365).describe()

```

```

Out[33]: 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

```

The data looks reasonable there is no outliers in DAYS_BIRTH.

```
In [0]: 1 plt.hist(Df_application_train['DAYS_BIRTH'] / 365, edgecolor = 'b', bir
2 plt.title('Age of Client in years'): plt.xlabel('Age'): plt.ylabel('Cou
```



People of age group between 35-40 years takes maximum no of loans.

```
In [0]: 1 (Df_application_train['DAYS_EMPLOYED']).describe()
```

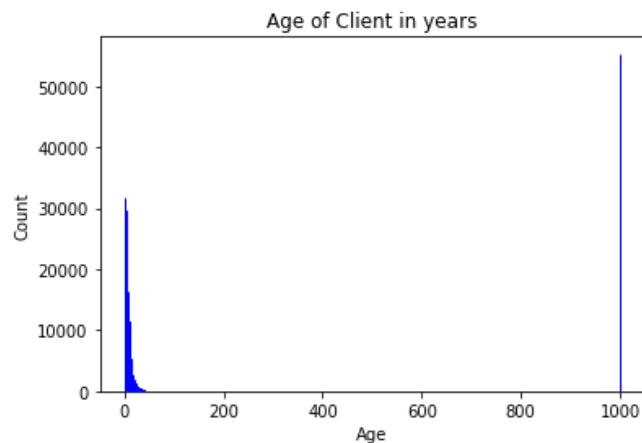
```
Out[21]: count    307511.000000
mean      63815.045904
std       141275.766519
min       -17912.000000
25%       -2760.000000
50%       -1213.000000
75%       -289.000000
max       365243.000000
Name: DAYS_EMPLOYED, dtype: float64
```

```
In [0]: 1 (Df_application_train['DAYS_EMPLOYED']).max()/365
```

```
Out[23]: 1000.6657534246575
```

365243 is max value for DAYS_EMPLOYED and after dividing it by 365 it gives 1000.66 which is an outlier because how can a employ works for 1000 years.

```
In [0]: 1 plt.hist(abs(Df_application_train['DAYS_EMPLOYED']) / 365, edgecolor =
2 plt.title('Age of Client in years'): plt.xlabel('Age'): plt.ylabel('Cou
```



The person with age 1000 is an outlier

```
In [0]: 1 Df_application_train['CODE GENDER'].unique()
```

```
Out[44]: array(['M', 'F', 'XNA'], dtype=object)
```

XNA is not a gender we should remove it

```
In [0]: 1 Df_application_train['NAME INCOME TYPE'].unique()
```

```
Out[47]: array(['Working', 'State servant', 'Commercial associate', 'Pensioner',
               'Unemployed', 'Student', 'Businessman', 'Maternity leave'],
              dtype=object)
```

Income type can never be Maternity leave so this should be removed.

```
In [0]: 1 d = [data for data in Df_application_train.columns if 'FLAG_DOC' in data]
        2
        3 Df_application_train[d]
```

```
Out[59]:
```

	FLAG_DOCUMENT_2	FLAG_DOCUMENT_3	FLAG_DOCUMENT_4	FLAG_DOCUMENT_5	FLAG_DOCUMENT_6
--	-----------------	-----------------	-----------------	-----------------	-----------------

0	0	1	0	0
1	0	1	0	0
2	0	0	0	0
3	0	1	0	0
4	0	0	0	0
...
307506	0	0	0	0
307507	0	1	0	0
307508	0	1	0	0
307509	0	1	0	0
307510	0	1	0	0

307511 rows × 20 columns

Checking no of 1's in all the flag documents.

```
In [0]: 1 for i in d:
2         a=sum(Df_application_train[i]==1)
3         print(i, "contain no of 1's = ", (a*100/307511), " no of 0's = ", (307511-
FLAG_DOCUMENT_2 contain no of 1's = 0.004227491049100683 no of 0's = 99.99577250895089
FLAG_DOCUMENT_3 contain no of 1's = 71.00233812774177 no of 0's = 28.997661872258227
FLAG_DOCUMENT_4 contain no of 1's = 0.008129790479039774 no of 0's = 99.99187020952095
FLAG_DOCUMENT_5 contain no of 1's = 1.5114906458630748 no of 0's = 98.48850935413692
FLAG_DOCUMENT_6 contain no of 1's = 8.80553866365756 no of 0's = 91.19446133634244
FLAG_DOCUMENT_7 contain no of 1's = 0.019186305530533868 no of 0's = 99.98081369446946
FLAG_DOCUMENT_8 contain no of 1's = 8.137595077899652 no of 0's = 91.86240492210035
FLAG_DOCUMENT_9 contain no of 1's = 0.389579559755586 no of 0's = 99.61042044024441
FLAG_DOCUMENT_10 contain no of 1's = 0.0022763413341311367 no of 0's = 99.99772365866586
FLAG_DOCUMENT_11 contain no of 1's = 0.39120551785139396 no of 0's = 99.60879448214861
FLAG_DOCUMENT_12 contain no of 1's = 0.000650383238323182 no of 0's = 99.99934961676168
FLAG_DOCUMENT_13 contain no of 1's = 0.3525077151711646 no of 0's = 99.64749228482883
FLAG_DOCUMENT_14 contain no of 1's = 0.2936480321029166 no of 0's = 99.70635196789708
FLAG_DOCUMENT_15 contain no of 1's = 0.12097128232811184 no of 0's = 99.87902871767189
FLAG_DOCUMENT_16 contain no of 1's = 0.9928100133003372 no of 0's = 99.00718998669966
FLAG_DOCUMENT_17 contain no of 1's = 0.02666571277125046 no of 0's = 99.97333428722875
FLAG_DOCUMENT_18 contain no of 1's = 0.8129790479039775 no of 0's = 99.18702095209602
FLAG_DOCUMENT_19 contain no of 1's = 0.059510066306571144 no of 0's = 99.94048993369343
FLAG_DOCUMENT_20 contain no of 1's = 0.05072989258920819 no of 0's = 99.9492701074108
FLAG_DOCUMENT_21 contain no of 1's = 0.03349473677364387 no of 0's = 99.96650526322635
```

All documents contains very less no of 1's except FLAG_DOCUMENT_3 so we can remove them as they cannot give much information.

```
In [0]: 1 sum(Df_application_train['AMT_GOODS_PRICE'].isnull())# no of null value
```

Out[75]: 278

```

In [5]: 1 #https://www.kaggle.com/aantonova/797-lgbm-and-bayesian-optimization
2 df = reduce_mem_usage(pd.read_csv('application_train.csv'))
3 test_df = reduce_mem_usage(pd.read_csv('application_test.csv'))
4 df = df.append(test_df).reset_index()
5 df = df[df['CODE_GENDER'] != 'XNA']
6 df = df[df['AMT_GOODS_PRICE'].notnull()]
7 df = df[df['NAME_INCOME_TYPE'] != 'Maternity leave']
8 df = df[df['DAYS_LAST_PHONE_CHANGE'].notnull()]
9 df['DAYS_EMPLOYED'].replace(365243, np.nan, inplace=True)
10 #If ratio of days birth to days employed is more that means person is l
11 df['DAYS_BIRTH / DAYS_EMPLOYED'] = df['DAYS_BIRTH'] / df['DAYS_EMPLOYED
12
13 df['AMT_CREDIT / AMT_INCOME_TOTAL'] = df['AMT_CREDIT'] / df['AMT_INCOME
14 df['CNT_FAM_MEMBERS / AMT_INCOME_TOTAL'] = df['CNT_FAM_MEMBERS'] / df['
15 df['AMT_ANNUITY / AMT_INCOME_TOTAL'] = df['AMT_ANNUITY'] / (1 + df['AMT
16
17 df['AMT_CREDIT / AMT_ANNUITY'] = df['AMT_CREDIT'] / df['AMT_ANNUITY']
18 df['ANNUITY_LENGTH / DAYS_EMPLOYED'] = df['AMT_CREDIT / AMT_ANNUITY'] /
19 df['CNT_CHILDREN / CNT_FAM_MEMBERS'] = df['CNT_CHILDREN'] / df['CNT_FAM
20
21 df['AMT_CREDIT / AMT_GOODS_PRICE'] = df['AMT_CREDIT'] / df['AMT_GOODS_F
22 df['AMT_CREDIT / AMT_GOODS_PRICE'] = df['AMT_CREDIT'] - df['AMT_GOODS_F
23 df['DAYS_REGISTRATION / DAYS_ID_PUBLISH'] = df['DAYS_REGISTRATION'] / c
24 df['DAYS_BIRTH / DAYS_REGISTRATION'] = df['DAYS_BIRTH'] / df['DAYS_REGI
25 df['DOC_SUM'] = df['FLAG_DOCUMENT_2'] + df['FLAG_DOCUMENT_3'] + df['FLA
26
27 df['NAN_AMT_ANNUITY'] = 1.0*np.isnan(df['AMT_ANNUITY'])
28 df['AGE_FINISH'] = df['DAYS_BIRTH']*(-1.0/365) + (df['AMT_CREDIT']/df['
29
30 d = [data for data in df.columns if 'FLAG_DOC' in data]
31 l = [data for data in df.columns if ('FLAG_' in data) & ('FLAG_DOC' not
32 df['NEW_DOC_IND_KURT'] = df[d].kurtosis(axis=1)
33 df['NEW_LIVE_IND_SUM'] = df[l].sum(axis=1)
34 df['AMT_INCOME_TOTAL / CNT_CHILDREN'] = df['AMT_INCOME_TOTAL'] / (1 + c
35 org_type = df[['AMT_INCOME_TOTAL', 'ORGANIZATION_TYPE']].groupby('ORGAN
36 df['NEW_INC_BY_ORG'] = df['ORGANIZATION_TYPE'].map(org_type)
37 df['AMT_INCOME_TOTAL / NEW_INC_BY_ORG'] = df['AMT_INCOME_TOTAL'] / df['
38 df['OWN_CAR_AGE / DAYS_BIRTH'] = df['OWN_CAR_AGE'] / df['DAYS_BIRTH']
39 df['OWN_CAR_AGE / DAYS_EMPLOYED'] = df['OWN_CAR_AGE'] / df['DAYS_EMPLOY
40 df['DAYS_LAST_PHONE_CHANGE / DAYS_BIRTH'] = df['DAYS_LAST_PHONE_CHANGE'
41 df['DAYS_LAST_PHONE_CHANGE / DAYS_EMPLOYED'] = df['DAYS_LAST_PHONE_CHAN
42
43 for f in ['CODE_GENDER', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY']:
44     df[f], c = pd.factorize(df[f])
45
46 cat_col = [i for i in df.columns if df[i].dtype == 'object']
47 df = pd.get_dummies(df, columns= cat_col)
48
49 df= df.drop(['FLAG_DOCUMENT_2','FLAG_DOCUMENT_4','FLAG_DOCUMENT_5','FLA
50     'FLAG_DOCUMENT_8','FLAG_DOCUMENT_9','FLAG_DOCUMENT_10', 'FLAG_DOCUM
51     'FLAG_DOCUMENT_14','FLAG_DOCUMENT_15','FLAG_DOCUMENT_16','FLAG_DOCU
52     'FLAG_DOCUMENT_20','FLAG_DOCUMENT_21'],axis=1)
53 print(df.shape)
54 del test_df
55 gc.collect()
56

```

Memory usage of dataframe is 286.23 MB
 Memory usage after optimization is: 92.38 MB
 Decreased by 67.7%
 Memory usage of dataframe is 45.00 MB
 Memory usage after optimization is: 14.60 MB
 Decreased by 67.6%
 (355967, 244)

Out[5]: 0

ONE_HOT Encoding of categorical features

```
In [0]: 1 #one hot encode the categorical data
        2 def one_hot(df):
        3     original_col = list(df.columns)
        4     c = [c for c in df.columns if df[c].dtype == 'object']#find categor
        5     df = pd.get_dummies(df, columns= c, dummy_na= True)#one_hot_encode
        6     new_col = [c for c in df.columns if c not in original_col]
        7     return df. new_col
```

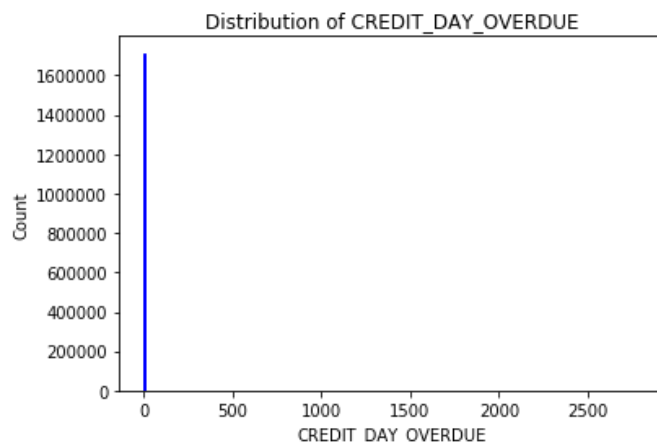
Bureau

```
In [0]: 1 b=pd.read_csv('bureau.csv')
```

```
In [0]: 1 b.columns
```

```
Out[11]: Index(['SK_ID_CURR', 'SK_ID_BUREAU', 'CREDIT_ACTIVE', 'CREDIT_CURRENCY',
               'DAYS_CREDIT', 'CREDIT_DAY_OVERDUE', 'DAYS_CREDIT_ENDDATE',
               'DAYS_ENDDATE_FACT', 'AMT_CREDIT_MAX_OVERDUE', 'CNT_CREDIT_PROLONG',
               'AMT_CREDIT_SUM', 'AMT_CREDIT_SUM_DEBT', 'AMT_CREDIT_SUM_LIMIT',
               'AMT_CREDIT_SUM_OVERDUE', 'CREDIT_TYPE', 'DAYS_CREDIT_UPDATE',
               'AMT_ANNUITY'],
              dtype='object')
```

```
In [0]: 1 plt.hist(b['CREDIT_DAY_OVERDUE'], edgecolor = 'b', bins = 1000)
        2 plt.title('Distribution of CREDIT DAY OVERDUE'): plt.xlabel('CREDIT DAY
```



Observation: Most of the credits have close to 0 days overdue

```
In [0]: 1 b.shape[0] #total no of entries in table
```

```
Out[20]: 1716428
```

```
In [0]: 1 b[b['CREDIT_DAY_OVERDUE']>0].shape[0] #entries which contain non zeroes
```

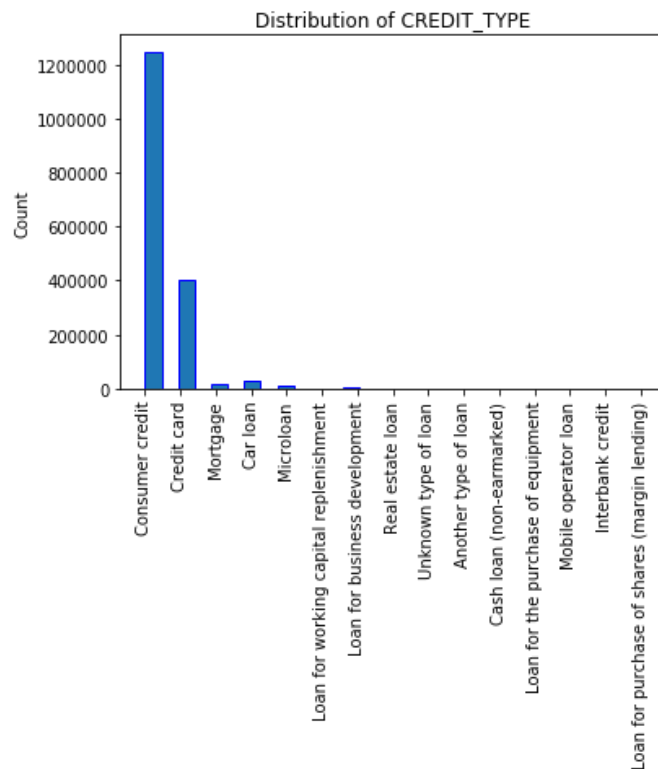
```
Out[19]: 4217
```

```
In [0]: 1 v=b['CREDIT_DAY_OVERDUE'].values
2 v=np.sort(v)
3 for i in range(99.,100,1):
4     print(i,'percentile value is',v[int(len(v)*float(i)/100)])
5     print(100,'percentile value is'.v[-1])

90 percentile value is 0
91 percentile value is 0
92 percentile value is 0
93 percentile value is 0
94 percentile value is 0
95 percentile value is 0
96 percentile value is 0
97 percentile value is 0
98 percentile value is 0
99 percentile value is 0
100 percentile value is 2792
```

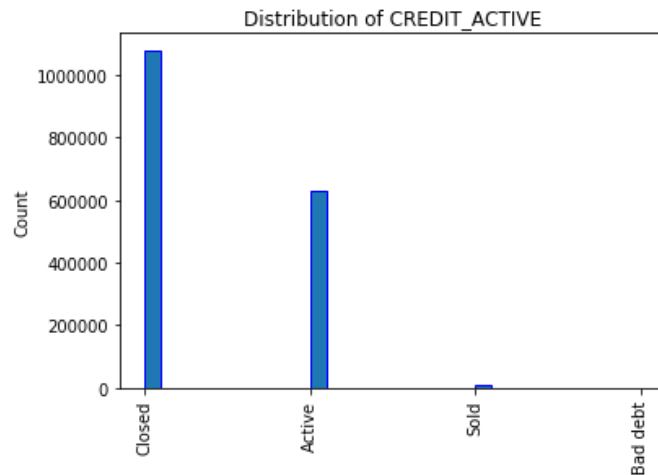
99% of values contain zero value.that means there is no overdue 99% of times.

```
In [0]: 1 plt.hist(b['CREDIT_TYPE'], edgecolor = 'b',stacked=0, bins = 30)
2 plt.title('Distribution of CREDIT_TYPE')
3 plt.xticks(rotation='vertical')
4 a=plt.ylabel('Count')
```



Observation: Consumer credit and Credit card are mostly registered credit in credit bureau.

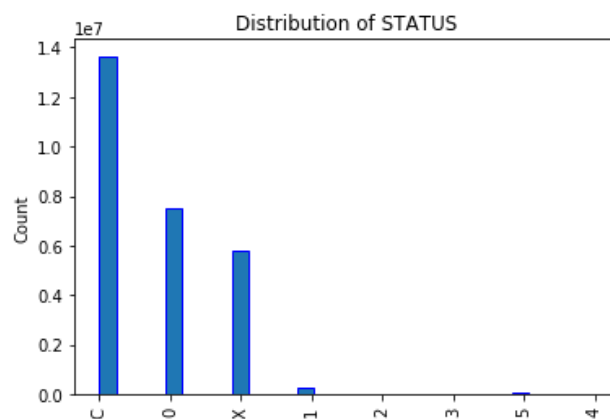
```
In [0]: 1
2 plt.hist(b['CREDIT_ACTIVE'], edgecolor = 'b',stacked=0, bins = 30)
3 plt.title('Distribution of CREDIT_ACTIVE')
4 plt.xticks(rotation='vertical')
5 a=plt.ylabel('Count')
```



Observation: Most of the credit registered are in status closed. Sold and bad debt are very few.

```
In [0]: 1 bb=pd.read_csv('/content/bureau_balance.csv')
```

```
In [0]: 1 plt.hist(bb['STATUS'], edgecolor = 'b',stacked=0, bins = 30)#Status of
2 #(active, closed, DPD0-30,à [C means closed, X means status unknown, 0
3 #1 means maximal did during month between 1-30, 2 means DPD 31-60,à 5 n
4 plt.title('Distribution of STATUS')
5 plt.xticks(rotation='vertical')
6 a=plt.ylabel('Count')
```



Observation: Most of the times credit bureau status is closed for a month

Features extraction in bureau_balance and bureau


```

In [0]: 1 #function takes the bureau and related dataframe and add aggregated fea
2 #also merge it with bureau_balance dataframe
3 #https://www.kaggle.com/jsaguiar/lightgbm-7th-place-solution
4 #https://www.linkedin.com/pulse/winning-9th-place-kaggles-biggest-compe
5 def b_new(buro):
6     bureau_balance = reduce_mem_usage(pd.read_csv('bureau_balance.csv'))
7     bureau_balance, col1 = one_hot(bureau_balance)
8     bureau, col2 = one_hot(buro)
9
10    bureau_balance_agg = {'MONTHS_BALANCE': ['size']}
11    for c in col1:
12        bureau_balance_agg[c] = ['mean']
13    aggregate = bureau_balance.groupby('SK_ID_BUREAU').agg(bureau_balanc
14    l=[]
15    for i in aggregate.columns.tolist():
16        l.append(i[0] + "_" + i[1].upper())
17
18    aggregate.columns = pd.Index(l)
19    bureau = bureau.join(aggregate, how='left', on='SK_ID_BUREAU')
20    bureau.drop(['SK_ID_BUREAU'], axis=1, inplace=True)
21    del aggregate
22    gc.collect()
23
24    numerical_agg = {'AMT_CREDIT_SUM_DEBT': ['mean', 'sum'], 'AMT_CREDIT
25        'DAYS_CREDIT': ['mean', 'var'], 'DAYS_CREDIT_UPDATE': ['mean'],
26        'DAYS_CREDIT_ENDDATE': ['mean'], 'CNT_CREDIT_PROLONG': ['sum'],
27        'AMT_CREDIT_SUM_LIMIT': ['mean', 'sum'], 'AMT_CREDIT_MAX_OVERDUE
28        'AMT_ANNUITY': ['max', 'mean'], 'AMT_CREDIT_SUM': ['mean', 'sum']
29    }
30    categorical_agg = {}
31    for c in col2:
32        categorical_agg[c] = ['mean']
33    for c in col1:
34        categorical_agg[c + "_MEAN"] = ['mean']
35
36    b1 = bureau.groupby('SK_ID_CURR').agg(**numerical_agg, **categoric
37    l=[]
38    for i in b1.columns.tolist():
39        l.append('BU_'+i[0]+'_'+i[1].upper())
40    b1.columns = pd.Index(l)
41
42    A = bureau[bureau['CREDIT_ACTIVE_Active'] == 1]
43    A_agg = A.groupby('SK_ID_CURR').agg(numerical_agg)
44    l=[]
45    for i in A_agg.columns.tolist():
46        l.append('A_'+i[0]+'_'+i[1].upper())
47    A_agg.columns = pd.Index(l)
48    b1 = b1.join(A_agg, how='left', on='SK_ID_CURR')
49    del A, A_agg
50    gc.collect()
51
52    C = bureau[bureau['CREDIT_ACTIVE_Closed'] == 1]
53    C_agg = C.groupby('SK_ID_CURR').agg(numerical_agg)
54    l=[]
55    for i in C_agg.columns.tolist():
56        l.append('C_'+i[0]+'_'+i[1].upper())
57    C_agg.columns = pd.Index(l)
58    b1 = b1.join(C_agg, how='left', on='SK_ID_CURR')
59    del C, C_agg
60    gc.collect()
61
62    print(b1.shape)
63    del bureau
64    gc.collect()
65    return(b1)
66

```

```

In [8]: 1
        2 b=pd.read_csv('bureau.csv')
        3
        4 new = b_new(reduce_mem_usage(b))
        5 df = df.join(new, how='left', on='SK_ID_CURR')
        6 del new
        7 gc.collect()

```

Memory usage of dataframe is 222.62 MB
 Memory usage after optimization is: 112.95 MB
 Decreased by 49.3%
 Memory usage of dataframe is 624.85 MB
 Memory usage after optimization is: 338.46 MB
 Decreased by 45.8%
 (305811, 89)

Out[8]: 0

Previous application

```

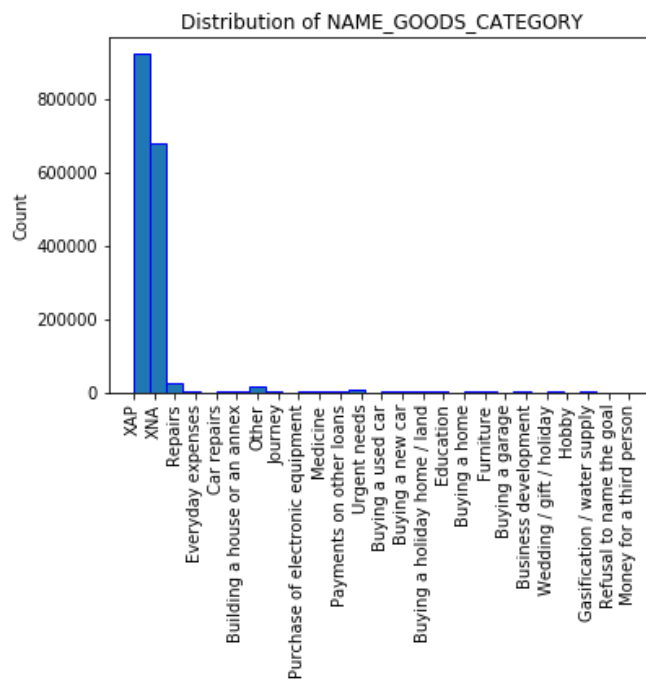
In [0]: 1 p=pd.read_csv('/content/previous_application.csv')

```

```

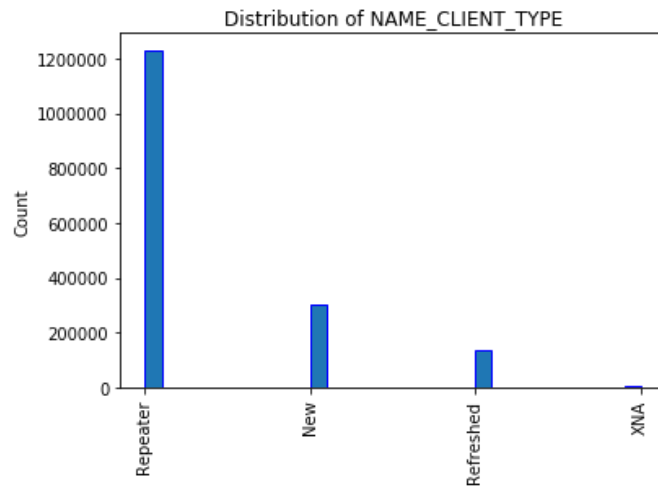
In [0]: 1
        2 plt.hist(p['NAME_CASH_LOAN_PURPOSE'], edgecolor = 'b',stacked=0, bins =
        3 plt.title('Distribution of NAME_GOODS_CATEGORY')
        4 plt.xticks(rotation='vertical')
        5 a=plt.vlabel('Count')

```



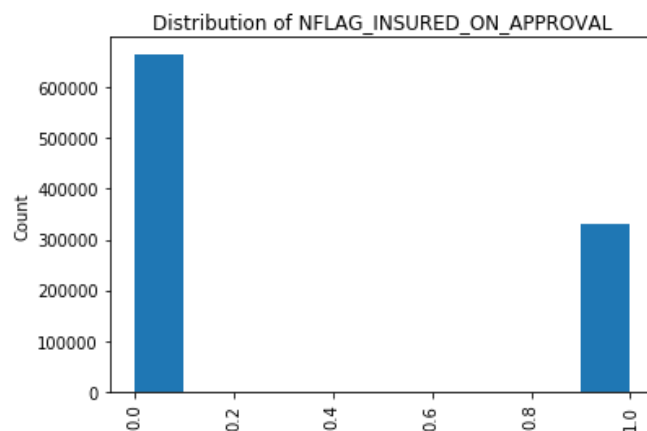
Observation: Most applicant takes cash loans with a purpose of XAP and XNA.

```
In [0]: 1 plt.hist(p['NAME_CLIENT_TYPE'], edgecolor = 'b', stacked=0, bins = 30)
2 plt.title('Distribution of NAME_CLIENT_TYPE')
3 plt.xticks(rotation='vertical')
4 a=plt.vlabel('Count')
```



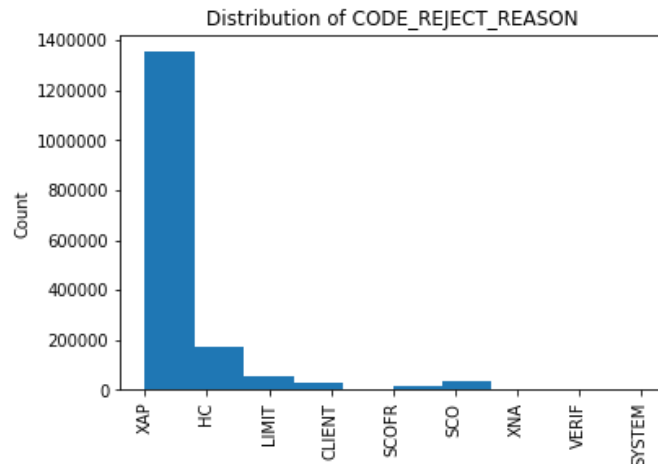
Observation: There are more no of repeaters when applying for the previous application.

```
In [0]: 1 plt.hist(p['NFLAG_INSURED_ON_APPROVAL'])#Did the client requested insurance
2 plt.title('Distribution of NFLAG_INSURED_ON_APPROVAL')
3 plt.xticks(rotation='vertical')
4 a=plt.vlabel('Count')
```



Observation: There are very less client who requested insurance during the previous application.

```
In [0]: 1 plt.hist(p['CODE_REJECT_REASON'])#Why was the previous application reje
2 plt.title('Distribution of CODE_REJECT_REASON')
3 plt.xticks(rotation='vertical')
4 a=plt.vlabel('Count')
```



Observation: Most of the times bank has given not applicable as the reason for rejecting the loan.

```
In [0]: 1 p['DAYS_TERMINATION'].describe()
```

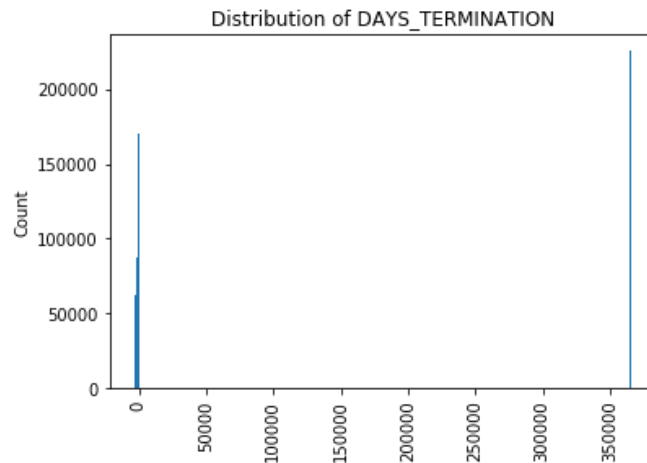
```
Out[64]: count    997149.000000
mean      81992.343838
std       153303.516729
min       -2874.000000
25%       -1270.000000
50%        -499.000000
75%        -44.000000
max       365243.000000
Name: DAYS_TERMINATION, dtype: float64
```

```
In [0]: 1 p[p['DAYS_TERMINATION']==365243]['DAYS_TERMINATION']
```

```
Out[85]: 1          365243.0
2          365243.0
17         365243.0
21         365243.0
34         365243.0
...
1669925     365243.0
1669945     365243.0
1669960     365243.0
1670192     365243.0
1670199     365243.0
Name: DAYS_TERMINATION, Length: 225913, dtype: float64
```

Observation: nearly 225913 values are outliers here in column DAYS_TERMINATION

```
In [0]: 1 plt.hist(p['DAYS_TERMINATION'],bins = 1000)#Relative to application dat
2 plt.title('Distribution of DAYS_TERMINATION')
3 plt.xticks(rotation='vertical')
4 a=plt.ylabel('Count')
```

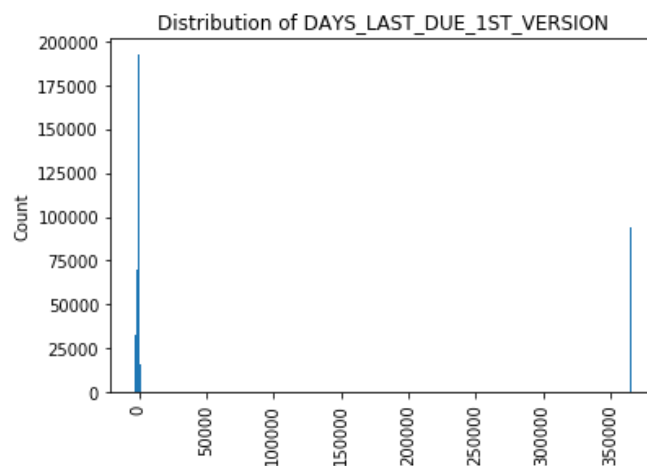


Observation: The min value is -2874 and max value is 365243. 365243 is equal to 1000 years hence it is an outlier.

```
In [0]: 1
```

```
In [0]: 1 print(p['DAYS_LAST_DUE_1ST_VERSION'].describe())
2 plt.hist(p['DAYS_LAST_DUE_1ST_VERSION'],bins = 1000)#Relative to applic
3 plt.title('Distribution of DAYS_LAST_DUE_1ST_VERSION')
4 plt.xticks(rotation='vertical')
5 a=plt.ylabel('Count')
```

```
count    997149.000000
mean      33767.774054
std      106857.034789
min       -2801.000000
25%       -1242.000000
50%        -361.000000
75%         129.000000
max       365243.000000
Name: DAYS_LAST_DUE_1ST_VERSION, dtype: float64
```



Observation: 365243 is an outlier here.

```
In [0]: 1 p['DAYS_FIRST_DRAWING'].describe()
```

```
Out[84]: count      997149.000000
         mean      342209.855039
         std       88916.115834
         min       -2922.000000
         25%       365243.000000
         50%       365243.000000
         75%       365243.000000
         max       365243.000000
         Name: DAYS_FIRST_DRAWING, dtype: float64
```

```
In [0]: 1 p[p['DAYS_FIRST_DRAWING']>0]['DAYS_FIRST_DRAWING']
```

```
Out[81]: 0          365243.0
         1          365243.0
         2          365243.0
         3          365243.0
         5          365243.0
         ...
         1670209     365243.0
         1670210     365243.0
         1670211     365243.0
         1670212     365243.0
         1670213     365243.0
         Name: DAYS_FIRST_DRAWING, Length: 934444, dtype: float64
```

```
In [0]: 1 p[p['DAYS_FIRST_DRAWING']==365243]['DAYS_FIRST_DRAWING']
```

```
Out[82]: 0          365243.0
         1          365243.0
         2          365243.0
         3          365243.0
         5          365243.0
         ...
         1670209     365243.0
         1670210     365243.0
         1670211     365243.0
         1670212     365243.0
         1670213     365243.0
         Name: DAYS_FIRST_DRAWING, Length: 934444, dtype: float64
```

```
In [0]: 1 p[p['DAYS_FIRST_DRAWING']<0]['DAYS_FIRST_DRAWING']
```

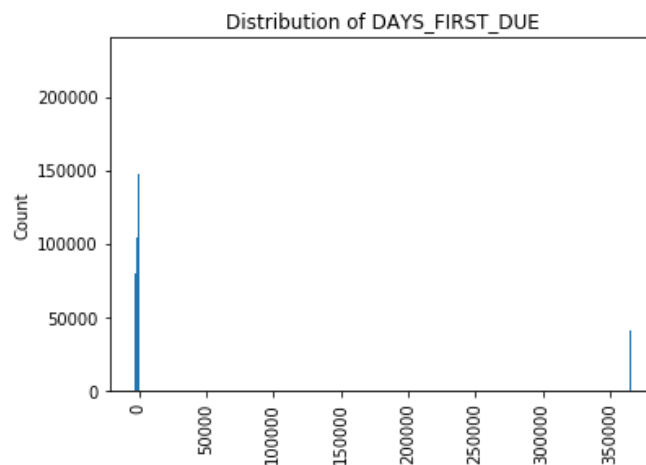
```
Out[83]: 17          -277.0
         34          -265.0
         82          -479.0
         93          -332.0
         242         -225.0
         ...
         1669788     -863.0
         1669830     -453.0
         1669833     -596.0
         1669960    -1083.0
         1670192     -474.0
         Name: DAYS_FIRST_DRAWING, Length: 62705, dtype: float64
```

All the values above 0 are 365243 and they all are outliers. Nearly 934444 values in DAYS_FIRST_DRAWING column are outliers except 62705 values which are less than 0. the min value here is -2922

```
In [0]: 1 p['DAYS_FIRST_DUE'].describe()
```

```
Out[87]: count      997149.000000
         mean      13826.269337
         std       72444.869708
         min       -2892.000000
         25%       -1628.000000
         50%        -831.000000
         75%       -411.000000
         max      365243.000000
         Name: DAYS_FIRST_DUE, dtype: float64
```

```
In [0]: 1 plt.hist(p['DAYS_FIRST_DUE'],bins = 1000)#Relative to application date
         2 plt.title('Distribution of DAYS_FIRST_DUE')
         3 plt.xticks(rotation='vertical')
         4 a=plt.ylabel('Count')
```



Observation: 365243 is outlier here.

```
In [0]: 1 print('There are',sum(p['DAYS_FIRST_DUE']<0),'reasonable value DAYS_FIR
         2 print("There are",sum(p['DAYS_FIRST_DUE']==365243),"outliers DAYS FIRST
```

```
There are 956504 reasonable value DAYS_FIRST_DUE
There are 40645 outliers DAYS_FIRST_DUE
```

```
In [0]: 1 print('There are',sum(p['DAYS_LAST_DUE']<0),'reasonable value in DAYS_L
         2 print("There are",sum(p['DAYS_LAST_DUE']==365243),"outliers in DAYS LAS
```

```
There are 785928 reasonable value in DAYS_LAST_DUE
There are 211221 outliers in DAYS_LAST_DUE
```

Feature engineering for Previous application

```

In [0]: 1 #function takes the previous_app and related dataframe and add aggregat
2 #https://www.kaggle.com/jsaguiar/lightgbm-7th-place-solution
3 #https://www.linkedin.com/pulse/winning-9th-place-kaggles-biggest-compe
4 def previous_app(prev):
5     p_app, p_app_cols = one_hot(prev)
6     p_app['DAYS_TERMINATION'].replace(365243, np.nan, inplace= True)
7     p_app['DAYS_LAST_DUE_1ST_VERSION'].replace(365243, np.nan, inplace=
8     p_app['DAYS_FIRST_DRAWING'].replace(365243, np.nan, inplace= True)
9     p_app['DAYS_FIRST_DUE'].replace(365243, np.nan, inplace= True)
10    p_app['DAYS_LAST_DUE'].replace(365243, np.nan, inplace= True)
11    p_app['APP_CREDIT_PERC'] = p_app['AMT_APPLICATION'] / p_app['AMT_CF
12    num_agg = {'AMT_DOWN_PAYMENT': ['max', 'mean'], 'DAYS_DECISION': ['n
13                'AMT_GOODS_PRICE': ['max', 'mean'], 'AMT_ANNUITY
14                'RATE_DOWN_PAYMENT': ['max', 'mean'], 'AMT_APPLICATION': ['max',
15                'AMT_CREDIT': ['max', 'mean'], 'APP_CREDIT_PERC': ['max', 'mean'
16
17    }
18    cat_agg = {}
19    for cat in p_app_cols:
20        cat_agg[cat] = ['mean']
21
22    P1 = p_app.groupby('SK_ID_CURR').agg(**num_agg, **cat_agg)
23    l=[]
24    for i in P1.columns.tolist():
25        l.append('PR_'+i[0]+'_'+i[1].upper())
26    P1.columns = pd.Index(l)
27    A = p_app[p_app['NAME_CONTRACT_STATUS_Approved'] == 1]
28    A_agg = A.groupby('SK_ID_CURR').agg(num_agg)
29    l=[]
30    for i in A_agg.columns.tolist():
31        l.append('APP_'+i[0]+'_'+i[1].upper())
32    A_agg.columns = pd.Index(l)
33    P1 = P1.join(A_agg, how='left', on='SK_ID_CURR')
34    R = p_app[p_app['NAME_CONTRACT_STATUS_Refused'] == 1]
35    R_agg = R.groupby('SK_ID_CURR').agg(num_agg)
36    l=[]
37    for i in R_agg.columns.tolist():
38        l.append('REF_'+i[0]+'_'+i[1].upper())
39
40    R_agg.columns = pd.Index(l)
41    P1 = P1.join(R_agg, how='left', on='SK_ID_CURR')
42    del R, R_agg, A, A_agg, p_app
43    gc.collect()
44    print(P1.shape)
45    return P1
46

```



```

In [11]: 1  # idea refered from https://www.kaggle.com/c/home-credit-default-risk
2  pr = pd.read_csv('previous_application.csv')
3
4  frame = previous_app(reduce_mem_usage(pr))
5  df = df.join(frame, how='left', on='SK_ID_CURR')
6  del frame
7  gc.collect()
8
9  frame = reduce_mem_usage(pr[pr.DAYS_DECISION >=-365].reset_index())
10 frame.drop(['index'],axis=1,inplace=True)
11 frame = previous_app(frame)
12 df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_year')
13 del frame
14 gc.collect()
15
16 frame = reduce_mem_usage(pr[pr.DAYS_DECISION >=-182].reset_index())
17 frame.drop(['index'],axis=1,inplace=True)
18 frame = previous_app(frame)
19 df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_halfyear')
20 del frame
21 gc.collect()
22
23 frame = reduce_mem_usage(pr[pr.DAYS_DECISION >=-92].reset_index())
24 frame.drop(['index'],axis=1,inplace=True)
25 frame = previous_app(frame)
26 df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_quarter')
27 del frame
28 gc.collect()
29
30 frame = reduce_mem_usage(pr[pr.DAYS_DECISION >=-31].reset_index())
31 frame.drop(['index'],axis=1,inplace=True)
32 frame = previous_app(frame)
33 df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_month')
34 del frame
35 gc.collect()
36
37 frame = reduce_mem_usage(pr[pr.DAYS_DECISION >=-8].reset_index())
38 frame.drop(['index'],axis=1,inplace=True)
39 frame = previous_app(frame)
40 df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_week')
41 del frame
42 gc.collect()

```

Memory usage of dataframe is 471.48 MB
Memory usage after optimization is: 309.01 MB
Decreased by 34.5%
(338857, 219)
Memory usage of dataframe is 110.53 MB
Memory usage after optimization is: 108.34 MB
Decreased by 2.0%
(195522, 216)
Memory usage of dataframe is 42.63 MB
Memory usage after optimization is: 41.78 MB
Decreased by 2.0%
(103914, 213)
Memory usage of dataframe is 17.44 MB
Memory usage after optimization is: 17.01 MB
Decreased by 2.5%
(52584, 210)
Memory usage of dataframe is 7.42 MB
Memory usage after optimization is: 7.23 MB
Decreased by 2.5%
(23184, 208)
Memory usage of dataframe is 1.92 MB
Memory usage after optimization is: 1.88 MB
Decreased by 2.5%
(6135, 206)

Out[11]: 0

POS_CASH_balance

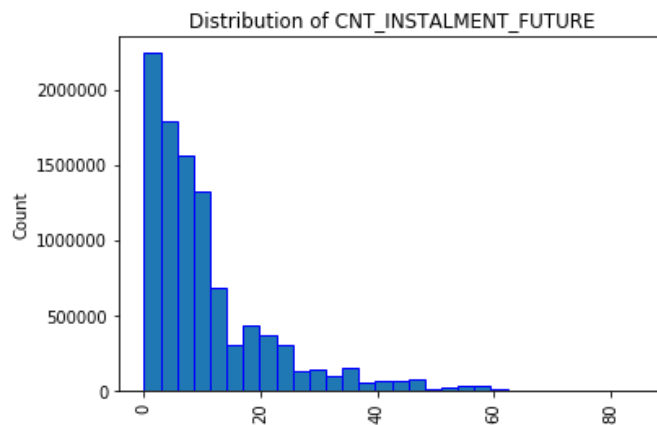
```
In [0]: 1 pc=pd.read_csv('/content/POS_CASH_balance.csv')
        2 pc
```

```
Out[51]:
```

	SK_ID_PREV	SK_ID_CURR	MONTHS_BALANCE	CNT_INSTALMENT	CNT_INSTALMENT_FUTURE
0	1803195	182943	-31	48.0	45.0
1	1715348	367990	-33	36.0	35.0
2	1784872	397406	-32	12.0	9.0
3	1903291	269225	-35	48.0	42.0
4	2341044	334279	-35	36.0	35.0
...
10001353	2448283	226558	-20	6.0	0.0
10001354	1717234	141565	-19	12.0	0.0
10001355	1283126	315695	-21	10.0	0.0
10001356	1082516	450255	-22	12.0	0.0
10001357	1259607	174278	-52	16.0	0.0

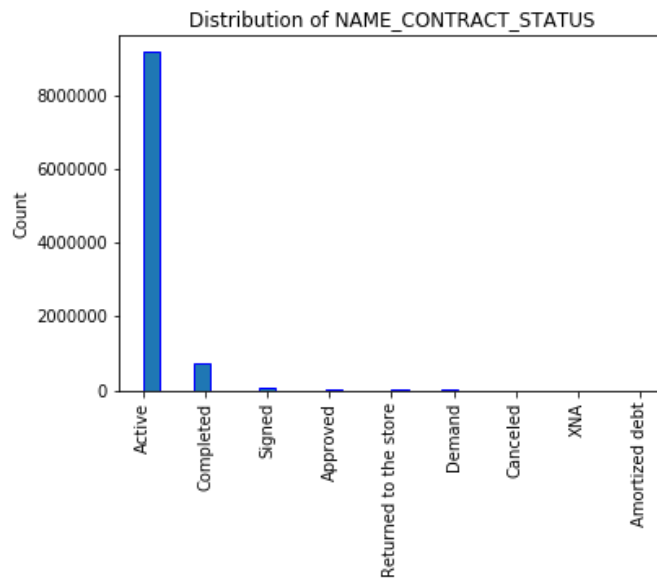
10001358 rows × 8 columns

```
In [0]: 1 plt.hist(pc['CNT_INSTALMENT_FUTURE'], edgecolor = 'b',stacked=0, bins =
        2 plt.title('Distribution of CNT_INSTALMENT_FUTURE')
        3 plt.xticks(rotation='vertical')
        4 a=plt.ylabel('Count')
```



Observation: installments left to pay mostly falls in the bracket of 0-20.

```
In [0]: 1 plt.hist(pc['NAME_CONTRACT_STATUS'], edgecolor = 'b',stacked=0, bins =
2 plt.title('Distribution of NAME_CONTRACT_STATUS')
3 plt.xticks(rotation='vertical')
4 a=plt.vlabel('Count')
```



Observation: Most of the time name contract status is Active

Feature engineering on POS_CASH_balance

```
In [0]: 1 #function takes the posh_cash and related dataframe and add aggregated
2 #https://www.kaggle.com/jsaguiar/lightgbm-7th-place-solution
3 #https://www.linkedin.com/pulse/winning-9th-place-kaggles-biggest-compe
4 def pos_cash(pos):
5     pos_ch, col = one_hot(pos)
6     aggregate = {'SK_DPD_DEF': ['max', 'mean'], 'SK_DPD': ['max', 'mean'],
7                   'MONTHS_BALANCE': ['max', 'mean', 'size']}
8
9
10    }
11    for c in col:
12        aggregate[c] = ['mean']
13
14    A_agg = pos_ch.groupby('SK_ID_CURR').agg(aggregate)
15    l=[]
16    for i in A_agg.columns.tolist():
17        l.append('POS_'+i[0] + '_' + i[1].upper())
18    A_agg.columns = pd.Index(l)
19    A_agg['POS_no_COUNT'] = pos_ch.groupby('SK_ID_CURR').size()
20
21    value=pos_ch.groupby('SK_ID_CURR')['MONTHS_BALANCE'].min()
22    A_agg['MIN_VALUES']=pos_ch['SK_ID_CURR'].map(value)
23    A_agg['MULTIPLIER']=1.00-pos_ch['MONTHS_BALANCE']/A_agg['MIN_VALUES']
24    print(A_agg.shape)
25    del pos_ch
26    gc.collect()
27    return A_agg
```

```

In [13]: 1      # idea refered from https://www.kaggle.com/c/home-credit-default-r
2      c = reduce_mem_usage(pd.read_csv('POS_CASH_balance.csv'))
3      frame = pos_cash(c)
4      df = df.join(frame, how='left', on='SK_ID_CURR')
5      del frame
6      gc.collect()
7
8      frame = reduce_mem_usage(c[c.MONTHS_BALANCE >=-12].reset_index())
9      frame.drop(['index'],axis=1,inplace=True)
10     frame = pos_cash(frame)
11     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_year')
12     del frame
13     gc.collect()
14
15     frame = reduce_mem_usage(c[c.MONTHS_BALANCE >=-6].reset_index())
16     frame.drop(['index'],axis=1,inplace=True)
17     frame = pos_cash(frame)
18     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_halfyear')
19     del frame
20     gc.collect()
21
22     frame = reduce_mem_usage(c[c.MONTHS_BALANCE >=-3].reset_index())
23     frame.drop(['index'],axis=1,inplace=True)
24     frame = pos_cash(frame)
25     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_quarter')
26     del frame
27     gc.collect()
28
29     frame = reduce_mem_usage(c[c.MONTHS_BALANCE >=-1].reset_index())
30     frame.drop(['index'],axis=1,inplace=True)
31     frame = pos_cash(frame)
32     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_month')
33     del frame
34     gc.collect()
35

```

```

Memory usage of dataframe is 610.43 MB
Memory usage after optimization is: 238.45 MB
Decreased by 60.9%
(337252, 20)
Memory usage of dataframe is 73.51 MB
Memory usage after optimization is: 64.60 MB
Decreased by 12.1%
(240235, 19)
Memory usage of dataframe is 33.01 MB
Memory usage after optimization is: 29.00 MB
Decreased by 12.1%
(192077, 19)
Memory usage of dataframe is 14.10 MB
Memory usage after optimization is: 12.39 MB
Decreased by 12.1%
(161387, 19)
Memory usage of dataframe is 2.99 MB
Memory usage after optimization is: 2.62 MB
Decreased by 12.1%
(79744, 18)

```

Out[13]: 0

Feature engineering on installments payments

```
In [0]: 1 #function takes the install_payments and related dataframe and add aggregate
2 #https://www.kaggle.com/jsaguiar/lightgbm-7th-place-solution
3 #https://www.linkedin.com/pulse/winning-9th-place-kaggles-biggest-competition
4 def inst_pay(ins):
5     A, col = one_hot(ins)
6
7     aggregate = {'days_late_on_payment': ['max', 'mean', 'min', 'sum'], 'AMT_INSTALLMENT_VERSION': ['nunique'], 'AMT_INSTALLMENT': ['max', 'sum'],
8                 'amount_extra_paid': ['max', 'mean', 'min', 'sum'], 'DAYS_ENTRY_PAYMENT': ['max', 'mean', 'min', 'sum']}
9
10    }
11    for c in col:
12        aggregate[c] = ['mean']
13    A_agg = A.groupby('SK_ID_CURR').agg(aggregate)
14    l=[]
15    for i in A_agg.columns.tolist():
16        l.append('INS_'+i[0]+'_'+i[1].upper())
17    A_agg.columns = pd.Index(l)
18    A_agg['INSTAL_no_COUNT'] = A.groupby('SK_ID_CURR').size()
19    del A
20    gc.collect()
21    print(A_agg.shape)
22    return A_agg
```

```

In [15]: 1      # idea refered from https://www.kaggle.com/c/home-credit-default-r
2      i=pd.read_csv('installments_payments.csv')
3      i['DAYS_ENTRY_PAYMENT'].fillna(100, inplace=True)
4      i['AMT_PAYMENT'].fillna(0.0, inplace=True)
5      i['days_late_on_payment']=i.DAYS_ENTRY_PAYMENT-i.DAYS_INSTALMENT
6      i['amount_extra_paid']=i.AMT_PAYMENT-i.AMT_INSTALMENT
7      i = reduce_mem_usage(i)
8      frame = inst_pay(i)
9      df = df.join(frame, how='left', on='SK_ID_CURR')
10     del frame
11     gc.collect()
12
13     frame = reduce_mem_usage(i[i.DAYS_INSTALMENT >=-365].reset_index())
14     frame.drop(['index'],axis=1,inplace=True)
15     frame = inst_pay(frame)
16     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_year')
17     del frame
18     gc.collect()
19
20     frame = reduce_mem_usage(i[i.DAYS_INSTALMENT >=-182].reset_index())
21     frame.drop(['index'],axis=1,inplace=True)
22     frame = inst_pay(frame)
23     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_halfyear')
24     del frame
25     gc.collect()
26
27     frame = reduce_mem_usage(i[i.DAYS_INSTALMENT >=-92].reset_index())
28     frame.drop(['index'],axis=1,inplace=True)
29     frame = inst_pay(frame)
30     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_quarter')
31     del frame
32     gc.collect()
33
34     frame = reduce_mem_usage(i[i.DAYS_INSTALMENT >=-31].reset_index())
35     frame.drop(['index'],axis=1,inplace=True)
36     frame = inst_pay(frame)
37     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_month')
38     del frame
39     gc.collect()
40
41     frame = reduce_mem_usage(i[i.DAYS_INSTALMENT >=-8].reset_index())
42     frame.drop(['index'],axis=1,inplace=True)
43     frame = inst_pay(frame)
44     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_week')
45     del frame
46     gc.collect()
47

```

Memory usage of dataframe is 1038.01 MB
Memory usage after optimization is: 389.25 MB
Decreased by 62.5%
(339587, 20)

Memory usage of dataframe is 124.89 MB
Memory usage after optimization is: 111.74 MB
Decreased by 10.5%
(252761, 20)

Memory usage of dataframe is 60.96 MB
Memory usage after optimization is: 54.54 MB
Decreased by 10.5%
(211804, 20)

Memory usage of dataframe is 28.81 MB
Memory usage after optimization is: 25.78 MB
Decreased by 10.5%
(186521, 20)

Memory usage of dataframe is 8.38 MB
Memory usage after optimization is: 7.50 MB
Decreased by 10.5%
(153389, 20)

Memory usage of dataframe is 1.38 MB

credit_card_balance

```
In [0]: 1 cc=pd.read_csv('credit card balance.csv')
```

```
In [0]: 1 sum(cc['CNT_INSTALLMENT_MATURE_CUM']==0)
```

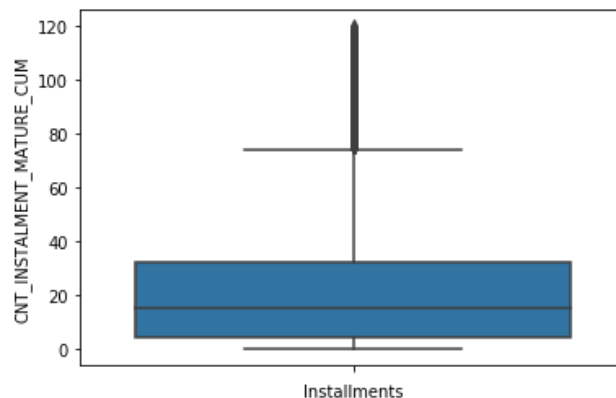
Out[55]: 551467

551467 times the paid installments are zero on the previous credit.

```
In [0]: 1 print(cc['CNT_INSTALLMENT_MATURE_CUM'].describe())
```

```
count    3.535076e+06
mean      2.082508e+01
std       2.005149e+01
min       0.000000e+00
25%       4.000000e+00
50%       1.500000e+01
75%       3.200000e+01
max       1.200000e+02
Name: CNT_INSTALLMENT_MATURE_CUM, dtype: float64
```

```
In [0]: 1 sns.boxplot(y='CNT_INSTALLMENT_MATURE_CUM', data=cc)
2 plt.xlabel('Installments')
3 plt.show()
```



Observation: 75% of the paid installments on the previous credit are below 40.

Feature engineering on credit credit balance

```
In [0]: 1 #function takes the credit_card and related dataframe and add aggregate
2 def credit_card(cre):
3     credit_card, col = one_hot(cre)
4     credit_card.drop(['SK_ID_PREV'], axis=1, inplace=True)
5     credit_card_agg = credit_card.groupby('SK_ID_CURR').agg(['max', 'me
6     l=[]
7     for i in credit_card_agg.columns.tolist():
8         l.append('CR_'+i[0]+'_'+i[1].upper())
9     credit_card_agg.columns = pd.Index(l)
10    # Count credit card lines
11    credit_card_agg['CR_no_COUNT'] = credit_card.groupby('SK_ID_CURR').
12    print(credit_card_agg.shape)
13    del credit_card
14    gc.collect()
15    return credit_card_agg
```

```

In [17]: 1      # idea refered from https://www.kaggle.com/c/home-credit-default-r
2      cr=pd.read_csv('credit_card_balance.csv')
3      frame = credit_card(reduce_mem_usage(cr))
4      df = df.join(frame, how='left', on='SK_ID_CURR')
5      del frame
6      gc.collect()
7
8      frame = reduce_mem_usage(cr[cr.MONTHS_BALANCE >=-12].reset_index())
9      frame.drop(['index'],axis=1,inplace=True)
10     frame = credit_card(frame)
11     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_year')
12     del frame
13     gc.collect()
14
15     frame = reduce_mem_usage(cr[cr.MONTHS_BALANCE >=-6].reset_index())
16     frame.drop(['index'],axis=1,inplace=True)
17     frame = credit_card(frame)
18     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_halfyear')
19     del frame
20     gc.collect()
21
22     frame = reduce_mem_usage(cr[cr.MONTHS_BALANCE >=-3].reset_index())
23     frame.drop(['index'],axis=1,inplace=True)
24     frame = credit_card(frame)
25     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_quarter')
26     del frame
27     gc.collect()
28
29     frame = reduce_mem_usage(cr[cr.MONTHS_BALANCE >=-1].reset_index())
30     frame.drop(['index'],axis=1,inplace=True)
31     frame = credit_card(frame)
32     df = df.join(frame, how='left', on='SK_ID_CURR',rsuffix='_month')
33     del frame
34     gc.collect()

```

Memory usage of dataframe is 673.88 MB
Memory usage after optimization is: 289.33 MB
Decreased by 57.1%
(103558, 113)
Memory usage of dataframe is 88.60 MB
Memory usage after optimization is: 84.53 MB
Decreased by 4.6%
(103558, 101)
Memory usage of dataframe is 46.35 MB
Memory usage after optimization is: 44.22 MB
Decreased by 4.6%
(103556, 101)
Memory usage of dataframe is 21.35 MB
Memory usage after optimization is: 20.37 MB
Decreased by 4.6%
(100677, 101)
Memory usage of dataframe is 5.17 MB
Memory usage after optimization is: 4.94 MB
Decreased by 4.6%
(61885, 101)

Out[17]: 0

Selecting top 300 features

```

In [18]: 1 df = reduce_mem_usage(df)
2

```

Memory usage of dataframe is 4718.39 MB
Memory usage after optimization is: 1996.12 MB
Decreased by 57.7%


```
In [0]: 1 import pickle
2 """
3 with open('new_data.pickle', 'wb') as handle:
4     pickle.dump(df, handle)
5 """
6 with open('new_data.pickle', 'rb') as handle:
7     df = pickle.load(handle)
```

```
In [0]: 1 train_df = df[df['TARGET'].notnull()]
2 test_df = df[df['TARGET'].isnull()]
```

```
In [0]: 1 from sklearn.feature_selection import SelectKBest
2 from sklearn.feature_selection import f_classif
3
4 Xsp = (train_df.loc[:, train_df.columns != 'TARGET'])
5 ysp = (train_df['TARGET'])
6
7
```

```
In [0]: 1 Xsp=np.nan to num(Xsp)
```

Splitting the dataset into train, cv and test.

```
In [0]: 1 from sklearn.model_selection import train_test_split
2
3 X_1, X_test, y_1, y_test = train_test_split(Xsp, ysp, test_size=0.10, s
4
5 # split the train data set into cross validation train and cross valida
6 X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.10, str
```

```
In [7]: 1 del df,Xsp,train_df
2 gc.collect()
```

Out[7]: 0

```
In [42]: 1 print(X_tr.shape, y_tr.shape)
2 print(X_cv.shape, y_cv.shape)
3 print(X_test.shape, y_test.shape)

(248850, 2336) (248850,)
(27650, 2336) (27650,)
(30723, 2336) (30723,)
```

```
In [0]: 1 S = SelectKBest(f_classif, k=300)
2 X_tr=S.fit_transform(X_tr, y_tr)
3
4
```

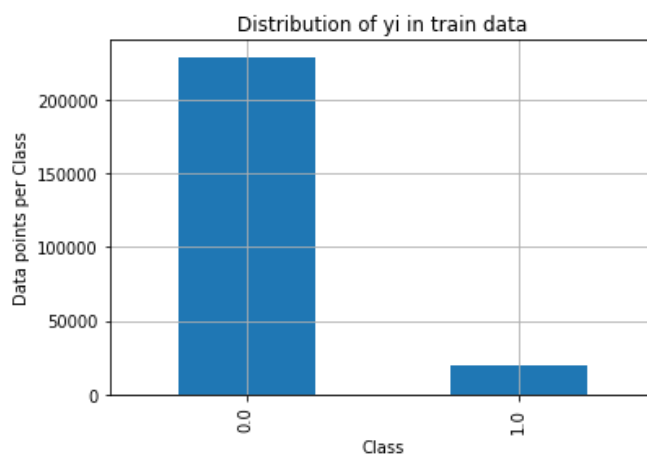
```
In [0]: 1 X_cv=S.transform(X_cv)
2 X_test=S.transform(X_test)
```

Checking the distribution of target variable into train,cv and test.

```

In [0]: 1 # it returns a dict, keys as class labels and values as the number of c
2 train_class_distribution = y_tr.value_counts()
3 test_class_distribution = y_test.value_counts()
4 cv_class_distribution = y_cv.value_counts()
5
6 my_colors = 'rbkymc'
7 train_class_distribution.plot(kind='bar')
8 plt.xlabel('Class')
9 plt.ylabel('Data points per Class')
10 plt.title('Distribution of yi in train data')
11 plt.grid()
12 plt.show()
13
14 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
15 # -(train_class_distribution.values): the minus sign will give us in de
16 sorted_yi = np.argsort(-train_class_distribution.values)
17 for i in sorted_yi:
18     print('Number of data points in class', i+1, ':', train_class_distrib
19
20
21 print('-'*80)
22 my_colors = 'rbkymc'
23 test_class_distribution.plot(kind='bar')
24 plt.xlabel('Class')
25 plt.ylabel('Data points per Class')
26 plt.title('Distribution of yi in test data')
27 plt.grid()
28 plt.show()
29
30 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
31 # -(train_class_distribution.values): the minus sign will give us in de
32 sorted_yi = np.argsort(-test_class_distribution.values)
33 for i in sorted_yi:
34     print('Number of data points in class', i+1, ':', test_class_distrib
35
36 print('-'*80)
37 my_colors = 'rbkymc'
38 cv_class_distribution.plot(kind='bar')
39 plt.xlabel('Class')
40 plt.ylabel('Data points per Class')
41 plt.title('Distribution of yi in cross validation data')
42 plt.grid()
43 plt.show()
44
45 # ref: argsort https://docs.scipy.org/doc/numpy/reference/generated/numpy.argsort.html
46 # -(train_class_distribution.values): the minus sign will give us in de
47 sorted_yi = np.argsort(-train_class_distribution.values)
48 for i in sorted_yi:
49     print('Number of data points in class', i+1, ':', cv_class_distribut
50

```



Logistic Regression

```
In [0]: 1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler()
3 scaler.fit(X_tr)
4 scaler_train = scaler.transform(X_tr)
5 scaler_cv = scaler.transform(X_cv)
6 scaler_test = scaler.transform(X_test)
```

```
In [0]: 1 def batch_predict(clf, data):
2     # roc_auc_score(y_true, y_score) the 2nd parameter should be probab
3     # not the predicted outputs
4
5     y_data_pred = []
6     tr_loop = data.shape[0] - data.shape[0]%1000
7     # consider you X_tr shape is 49041, then your cr_loop will be 49041
8     # in this for loop we will iterate until the last 1000 multiplier
9     for i in range(0, tr_loop, 1000):
10        y_data_pred.extend(clf.predict_proba(data[i:i+1000])[:,1])
11    # we will be predicting for the last data points
12    y_data_pred.extend(clf.predict_proba(data[tr_loop:])[:,1])
13
14    return y_data_pred
```

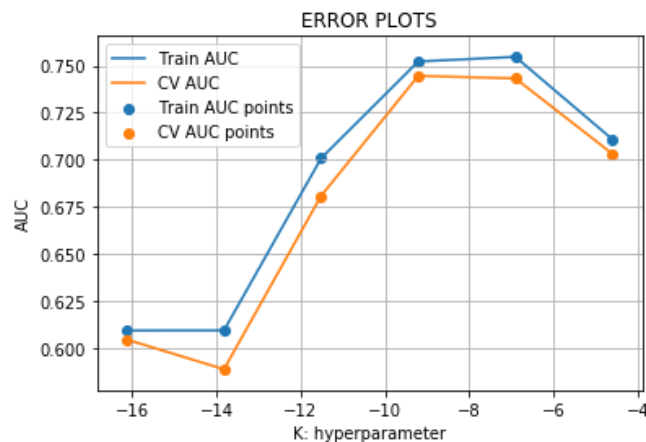
```
In [0]: 1 from sklearn.metrics import roc_curve, auc
2 from sklearn.metrics import roc_auc_score
3 from sklearn import linear_model
4 from tadm import tadm
```

```

In [0]: 1 train_auc = []
        2 cv_auc = []
        3 K = {'alpha': [10**-7,10**-6,10**-5,10**-4,10**-3,10**-2]}
        4 for i in tqdm(K['alpha']):
        5     neigh = linear_model.SGDClassifier(loss='log',alpha=i,penalty='l1')
        6     neigh.fit(scaler_train, y_tr)
        7
        8     y_train_pred = batch_predict(neigh, scaler_train)
        9     y_cv_pred = batch_predict(neigh, scaler_cv)
        10
        11     train_auc.append(roc_auc_score(y_tr,y_train_pred))
        12     cv_auc.append(roc_auc_score(y_cv, y_cv_pred))
        13 log_K=[]
        14 for l in K['alpha']:
        15     log_K.append(math.log(l))
        16 plt.plot(log_K, train_auc, label='Train AUC')
        17 plt.plot(log_K, cv_auc, label='CV AUC')
        18
        19 plt.scatter(log_K, train_auc, label='Train AUC points')
        20 plt.scatter(log_K, cv_auc, label='CV AUC points')
        21
        22 plt.legend()
        23 plt.xlabel("K: hyperparameter")
        24 plt.ylabel("AUC")
        25 plt.title("ERROR PLOTS")
        26 plt.grid()
        27 plt.show()

```

100%|██████████| 6/6 [05:52<00:00, 48.31s/it]



```

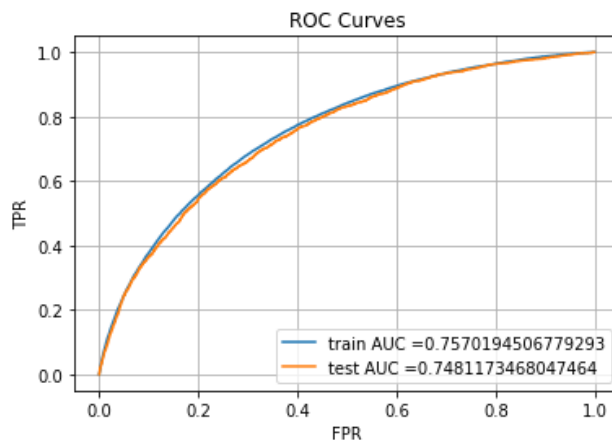
In [0]: 1 best_parl=10**-3

```

```

In [33]: 1 from sklearn.metrics import roc_curve, auc
          2
          3 neigh = linear_model.SGDClassifier(loss='log', alpha=best_par1, class_w
          4 neigh.fit(scaler_train, y_tr)
          5 # roc_auc_score(y_true, y_score) the 2nd parameter should be probabilit
          6 # not the predicted outputs
          7
          8 y_train_pred = batch_predict(neigh, scaler_train)
          9 y_test_pred = batch_predict(neigh, scaler_test)
         10
         11 train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, y_train_pred)
         12 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
         13
         14 plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, t
         15 plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_
         16 plt.legend()
         17 plt.xlabel("FPR")
         18 plt.ylabel("TPR")
         19 plt.title("ROC Curves")
         20 plt.grid()
         21 plt.show()

```



```

In [0]: 1 def predict(proba, threshold, fpr, tpr):
          2
          3     t = threshold[np.argmax(tpr*(1-fpr))]
          4
          5     # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is
          6
          7     # print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for
          8     predictions = []
          9     for i in proba:
         10         if i >= t:
         11             predictions.append(1)
         12         else:
         13             predictions.append(0)
         14     return predictions

```

Train

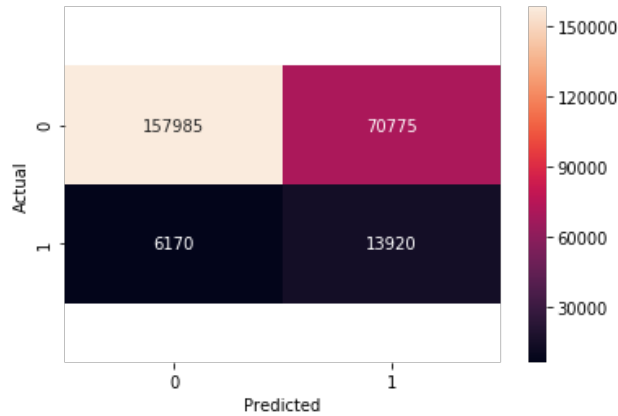
```

In [35]: 1 import seaborn as sns
          2 print("Train confusion matrix")
          3 ax=sns.heatmap(confusion_matrix(y_tr, predict(y_train_pred, tr_threshol
          4 bottom, top = ax.get_ylim()
          5 ax.set_ylim(bottom + 0.5, top - 0.5)
          6 plt.xlabel("Predicted")
          7 plt.ylabel("Actual")
          8

```

Train confusion matrix

Out[35]: Text(33.0, 0.5, 'Actual')



Test

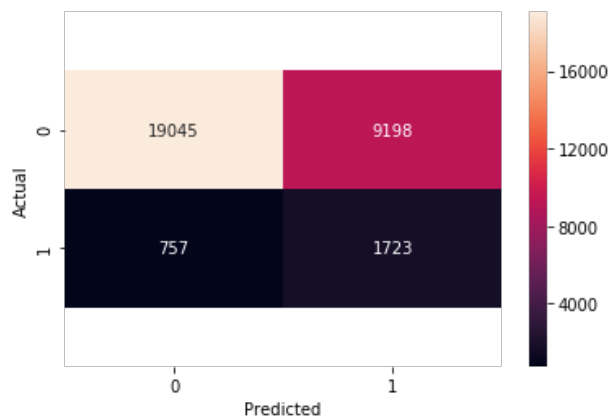
```

In [36]: 1 import seaborn as sns
          2 print("Train confusion matrix")
          3 ax=sns.heatmap(confusion_matrix(y_test, predict(y_test_pred, te_thresho
          4 bottom, top = ax.get_ylim()
          5 ax.set_ylim(bottom + 0.5, top - 0.5)
          6 plt.xlabel("Predicted")
          7 plt.ylabel("Actual")
          8

```

Train confusion matrix

Out[36]: Text(33.0, 0.5, 'Actual')



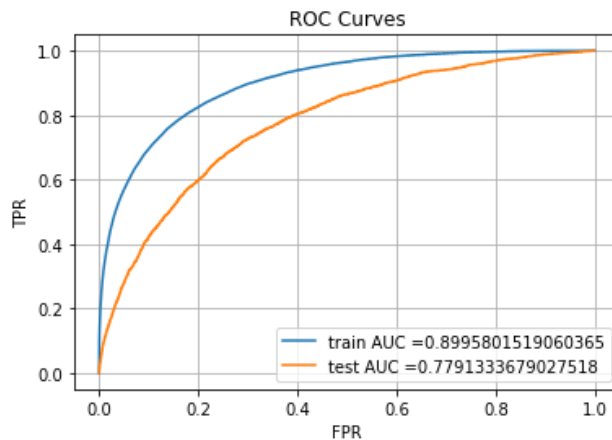
XGBoost

```
In [1]: 1 from sklearn.model_selection import StratifiedKFold
2
3 k=StratifiedKFold(n_splits=3, shuffle=True)
4 print(k)
StratifiedKFold(n_splits=3, random_state=None, shuffle=True)
```

```
In [21]: 1 from xgboost.sklearn import XGBClassifier
2 from sklearn.model_selection import StratifiedKFold
3 from sklearn.model_selection import RandomizedSearchCV
4 n_estimators = [700,800,900,1000]
5 learning_rate = [ 0.02,0.05,0.07]
6 max_depth = [6,7,8,9]
7
8 param = dict(learning_rate=learning_rate,n_estimators=n_estimators,max_
9 xgb = XGBClassifier(nthread=-1)
10 kfold = StratifiedKFold(n_splits=3, shuffle=True)
11 random_search_cv = RandomizedSearchCV(xgb, param, scoring="neg_log_loss
12 result = random_search_cv.fit(scaler_train , y_tr)
13
14 print("Best mean_test_score: ", result.best_score_ , "parameters", result
15 mean = result.cv_results_['mean_test_score']
16 std = result.cv_results_['std_test_score']
17 par = result.cv_results_['params']
18 for m, s, p in zip(mean, std, par):
19     print( 'mean_test_score',m,'std_test_score',s,'params',p)
20
```

```
Best mean_test_score: -0.2406262771923414 parameters {'n_estimators': 80
0, 'max_depth': 6, 'learning_rate': 0.05}
mean_test_score -0.24465255370463834 std_test_score 0.0008739492716115714
params {'n_estimators': 700, 'max_depth': 8, 'learning_rate': 0.05}
mean_test_score -0.24186662961652886 std_test_score 0.0007838895703288978
params {'n_estimators': 700, 'max_depth': 7, 'learning_rate': 0.05}
mean_test_score -0.24631985675540405 std_test_score 0.000790347957465164 p
arams {'n_estimators': 800, 'max_depth': 7, 'learning_rate': 0.07}
mean_test_score -0.24428878322120132 std_test_score 0.0007794105997974142
params {'n_estimators': 1000, 'max_depth': 6, 'learning_rate': 0.07}
mean_test_score -0.2448484095207295 std_test_score 0.0007974006162268901 p
arams {'n_estimators': 700, 'max_depth': 7, 'learning_rate': 0.07}
mean_test_score -0.24913924356126568 std_test_score 0.0006932261691023411
params {'n_estimators': 1000, 'max_depth': 8, 'learning_rate': 0.05}
mean_test_score -0.2677679252731944 std_test_score 0.0015705616987893593 p
arams {'n_estimators': 900, 'max_depth': 9, 'learning_rate': 0.07}
mean_test_score -0.2409831330670571 std_test_score 0.0009611824046316707 p
arams {'n_estimators': 900, 'max_depth': 6, 'learning_rate': 0.05}
mean_test_score -0.2406262771923414 std_test_score 0.000906651959100784 pa
rams {'n_estimators': 800, 'max_depth': 6, 'learning_rate': 0.05}
mean_test_score -0.24259295039010853 std_test_score 0.000758119901510508 p
arams {'n_estimators': 800, 'max_depth': 7, 'learning_rate': 0.05}
```

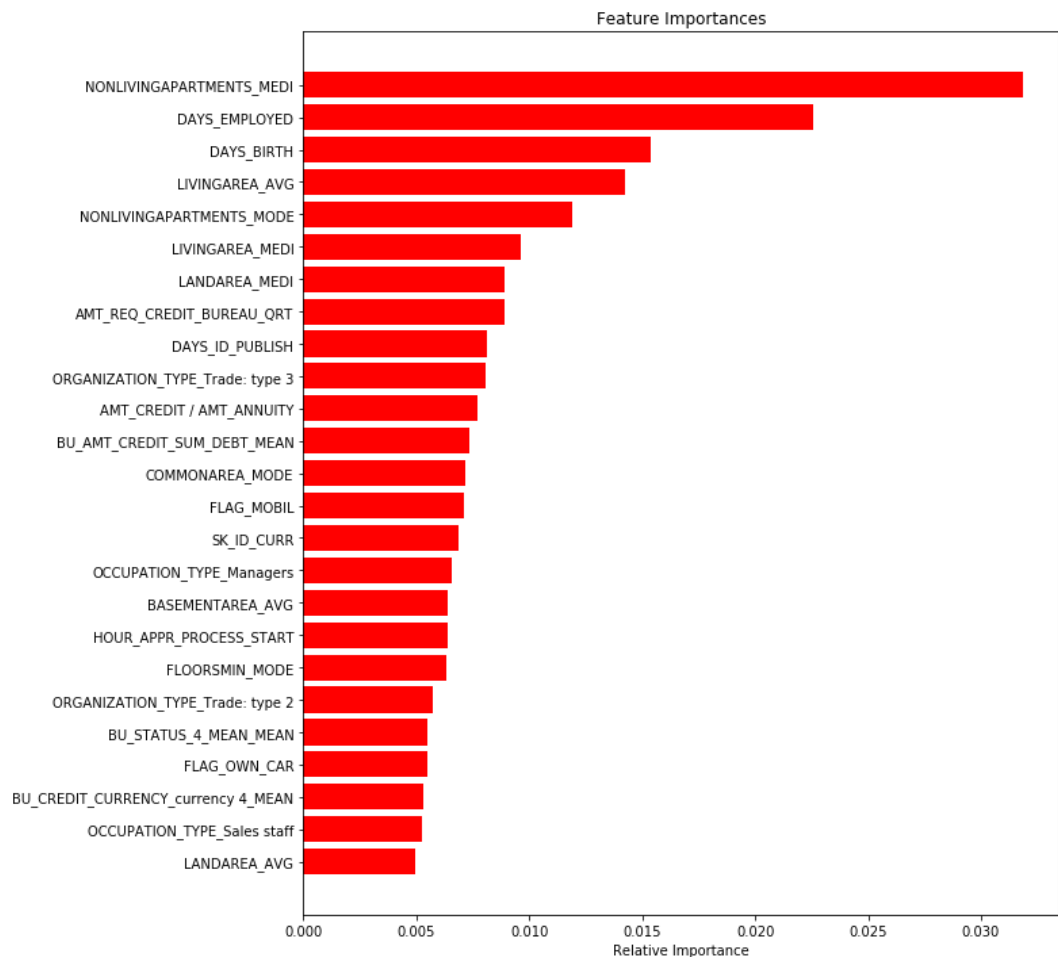
```
In [30]: 1 from sklearn.metrics import roc_curve, auc
2 from xgboost.sklearn import XGBClassifier
3 neigh = XGBClassifier(n_estimators=800,max_depth=6,learning_rate=0.05,
4                       )
5 neigh.fit(scaler_train, y_tr)
6 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability
7 # not the predicted outputs
8
9 y_train_pred = batch_predict(neigh, scaler_train)
10 y_test_pred = batch_predict(neigh, scaler_test)
11
12 train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, y_train_pred)
13 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
14
15 plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
16 plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
17 plt.legend()
18 plt.xlabel("FPR")
19 plt.ylabel("TPR")
20 plt.title("ROC Curves")
21 plt.grid()
22 plt.show()
```




```

In [19]: 1 import matplotlib.pyplot as plt
          2 features = train_df.columns
          3 importances = neigh.feature_importances_
          4 indices = (np.argsort(importances))[-25:]
          5 plt.figure(figsize=(10,12))
          6 plt.title('Feature Importances')
          7 plt.barh(range(len(indices)), importances[indices], color='r', align='c')
          8 plt.yticks(range(len(indices)), [features[i] for i in indices])
          9 plt.xlabel('Relative Importance')
         10 plt.show()

```

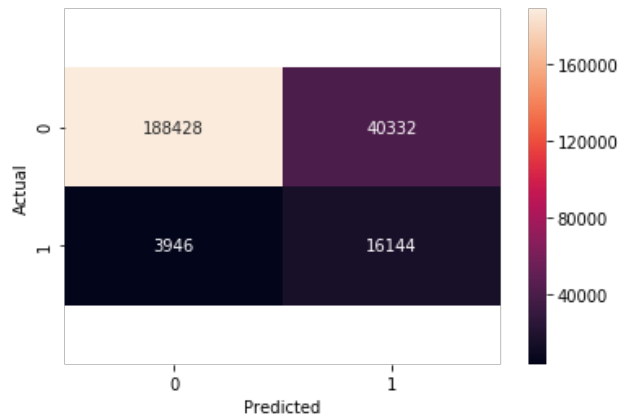


Train

```
In [20]: 1 import seaborn as sns
2 print("Train confusion matrix")
3 ax=sns.heatmap(confusion_matrix(y_tr, predict(y_train_pred, tr_threshol
4 bottom, top = ax.get_ylim()
5 ax.set_ylim(bottom + 0.5, top - 0.5)
6 plt.xlabel("Predicted")
7 plt.ylabel("Actual")
8
```

Train confusion matrix

Out[20]: Text(33.0, 0.5, 'Actual')

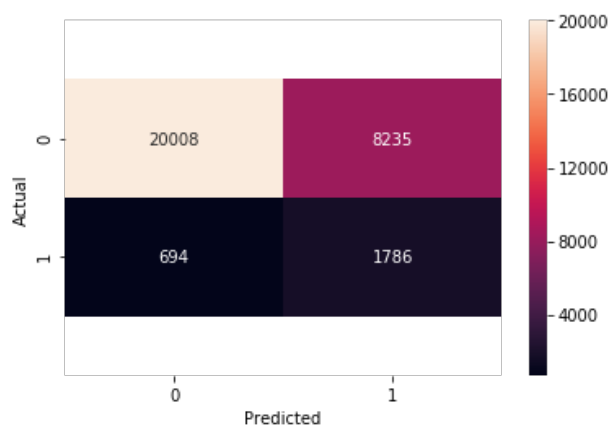


test

```
In [21]: 1 import seaborn as sns
2 print("Train confusion matrix")
3 ax=sns.heatmap(confusion_matrix(y_test, predict(y_test_pred, te_thresho
4 bottom, top = ax.get_ylim()
5 ax.set_ylim(bottom + 0.5, top - 0.5)
6 plt.xlabel("Predicted")
7 plt.ylabel("Actual")
8
```

Train confusion matrix

Out[21]: Text(33.0, 0.5, 'Actual')



Random Forest

```

In [23]: 1 from sklearn.ensemble import RandomForestClassifier
          2 from sklearn.model_selection import StratifiedKFold
          3 from sklearn.model_selection import RandomizedSearchCV
          4 n_estimators = [500,600,700,800]
          5 max_depth = [9,10,11,12]
          6
          7 param = dict(n_estimators=n_estimators,max_depth=max_depth)
          8 xgb = RandomForestClassifier(n_jobs=-1)
          9 kfold = StratifiedKFold(n_splits=3, shuffle=True)
         10 random_search_cv = RandomizedSearchCV(xgb, param, scoring="neg_log_loss")
         11 result = random_search_cv.fit(scaler_train , y_tr)
         12
         13 print("Best mean_test_score: ", result.best_score_ , "parameters", result.best_params_)
         14 mean = result.cv_results_['mean_test_score']
         15 std = result.cv_results_['std_test_score']
         16 par = result.cv_results_['params']
         17 for m, s, p in zip(mean, std, par):
         18     print( 'mean_test_score',m,'std_test_score',s,'params',p)
         19

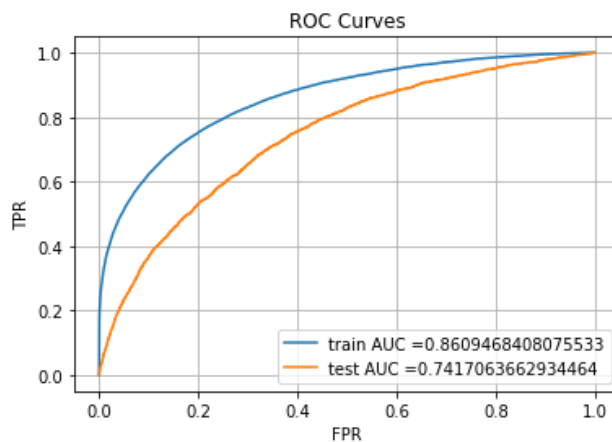
```

```

Best mean_test_score: -0.2536874293927612 parameters {'n_estimators': 800, 'max_depth': 12}
mean_test_score -0.254615364761944 std_test_score 0.0002823405537752798 params {'n_estimators': 500, 'max_depth': 11}
mean_test_score -0.25651790067081387 std_test_score 0.0001920862401660547 params {'n_estimators': 500, 'max_depth': 9}
mean_test_score -0.2565721240025001 std_test_score 0.0003010964807867514 params {'n_estimators': 600, 'max_depth': 9}
mean_test_score -0.255477673191798 std_test_score 0.0001977923024301139 params {'n_estimators': 800, 'max_depth': 10}
mean_test_score -0.25376309773251715 std_test_score 0.0002614216605694983 params {'n_estimators': 600, 'max_depth': 12}
mean_test_score -0.2536874293927612 std_test_score 0.000288791520100807 params {'n_estimators': 800, 'max_depth': 12}
mean_test_score -0.2554554413614863 std_test_score 0.00010584947643292523 params {'n_estimators': 700, 'max_depth': 10}
mean_test_score -0.25379525139531606 std_test_score 0.00021145300208143608 params {'n_estimators': 700, 'max_depth': 12}
mean_test_score -0.25456206385962066 std_test_score 0.00021551356505262402 params {'n_estimators': 600, 'max_depth': 11}
mean_test_score -0.2565865410647776 std_test_score 0.00031073006522073344 params {'n_estimators': 800, 'max_depth': 9}

```

```
In [24]: 1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import roc_curve, auc
3 neigh = RandomForestClassifier(n_estimators=800,max_depth=12)
4
5 neigh.fit(scaler_train, y_tr)
6 # roc_auc_score(y_true, y_score) the 2nd parameter should be probability
7 # not the predicted outputs
8
9 y_train_pred = batch_predict(neigh, scaler_train)
10 y_test_pred = batch_predict(neigh, scaler_test)
11
12 train_fpr, train_tpr, tr_thresholds = roc_curve(y_tr, y_train_pred)
13 test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_pred)
14
15 plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
16 plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
17 plt.legend()
18 plt.xlabel("FPR")
19 plt.ylabel("TPR")
20 plt.title("ROC Curves")
21 plt.grid()
22 plt.show()
```

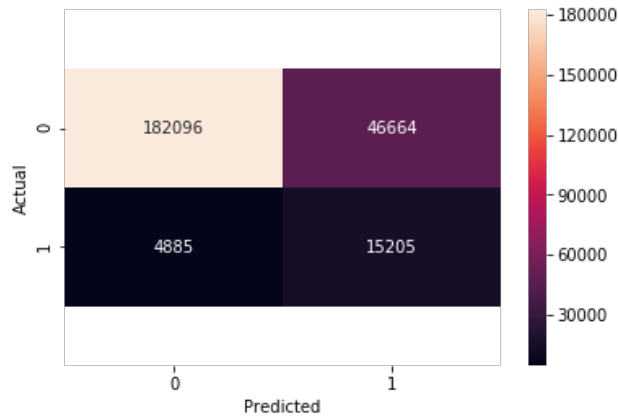


Train

```
In [25]: 1 import seaborn as sns
2 print("Train confusion matrix")
3 ax=sns.heatmap(confusion_matrix(y_tr, predict(y_train_pred, tr_threshol
4 bottom, top = ax.get_ylim()
5 ax.set_ylim(bottom + 0.5, top - 0.5)
6 plt.xlabel("Predicted")
7 plt.ylabel("Actual")
8
```

Train confusion matrix

Out[25]: Text(33.0, 0.5, 'Actual')

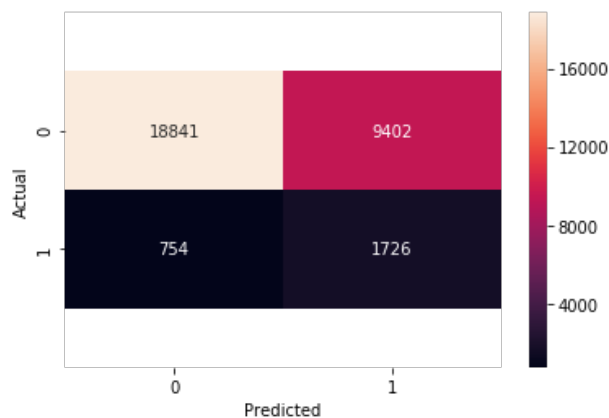


test

```
In [26]: 1 import seaborn as sns
2 print("Train confusion matrix")
3 ax=sns.heatmap(confusion_matrix(y_test, predict(y_test_pred, te_thresho
4 bottom, top = ax.get_ylim()
5 ax.set_ylim(bottom + 0.5, top - 0.5)
6 plt.xlabel("Predicted")
7 plt.ylabel("Actual")
8
```

Train confusion matrix

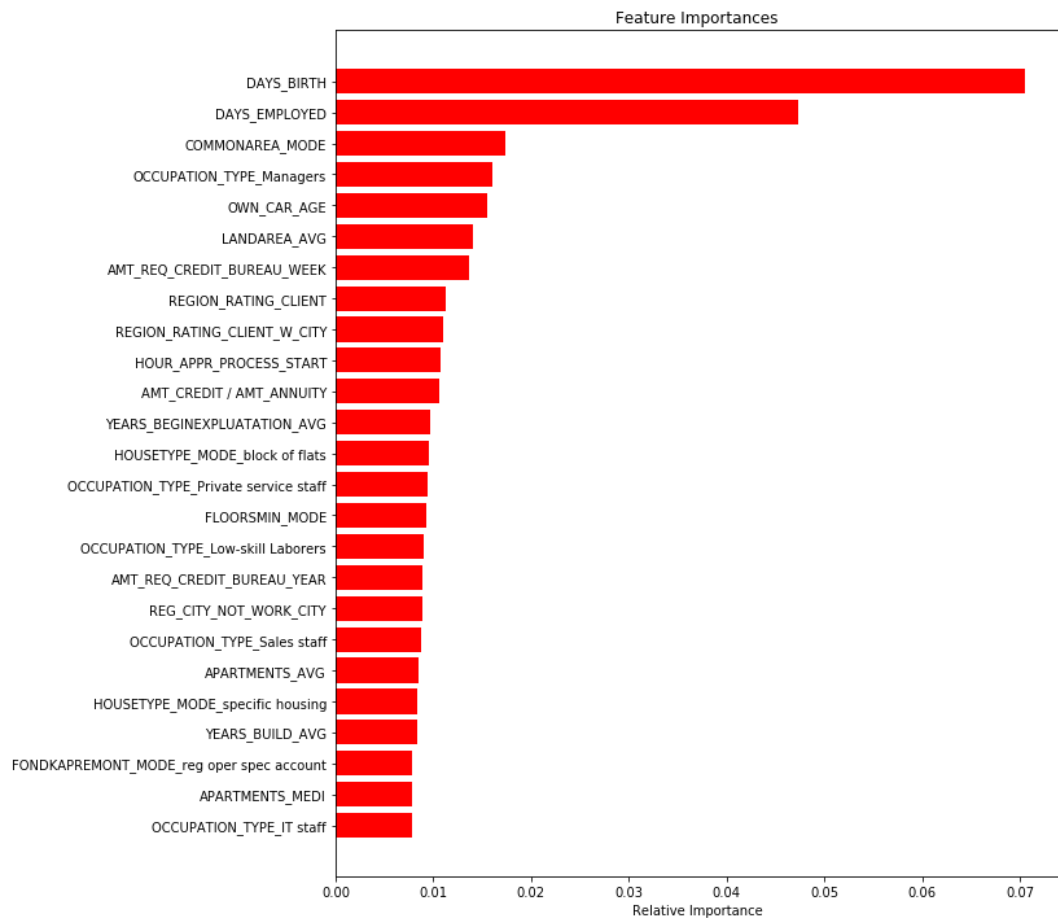
Out[26]: Text(33.0, 0.5, 'Actual')



```

In [29]: 1 import matplotlib.pyplot as plt
          2 features = train_df.columns
          3 importances = neigh.feature_importances_
          4 indices = (np.argsort(importances))[-25:]
          5 plt.figure(figsize=(10,12))
          6 plt.title('Feature Importances')
          7 plt.barh(range(len(indices)), importances[indices], color='r', align='c')
          8 plt.yticks(range(len(indices)), [features[i] for i in indices])
          9 plt.xlabel('Relative Importance')
         10 plt.show()

```



Hyper parameter tuning using BayesianOptimization

- Bayesian optimization methods build a probability model of the objective function to propose smarter choices for the next set of hyperparameters to evaluate. hence reduce the tuning time for hyperparameters.

```
In [22]: 1 !pip install bayesian-optimization
```

```
Collecting bayesian-optimization
  Downloading https://files.pythonhosted.org/packages/72/0c/173ac467d0a53e33e41b521e4ceba74a8ac7c7873d7b857a8fbdca88302d/bayesian-optimization-1.0.1.tar.gz (https://files.pythonhosted.org/packages/72/0c/173ac467d0a53e33e41b521e4ceba74a8ac7c7873d7b857a8fbdca88302d/bayesian-optimization-1.0.1.tar.gz)
Requirement already satisfied: numpy>=1.9.0 in /usr/local/lib/python3.6/dist-packages (from bayesian-optimization) (1.17.5)
Requirement already satisfied: scipy>=0.14.0 in /usr/local/lib/python3.6/dist-packages (from bayesian-optimization) (1.4.1)
Requirement already satisfied: scikit-learn>=0.18.0 in /usr/local/lib/python3.6/dist-packages (from bayesian-optimization) (0.22.1)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.6/dist-packages (from scikit-learn>=0.18.0->bayesian-optimization) (0.14.1)
Building wheels for collected packages: bayesian-optimization
  Building wheel for bayesian-optimization (setup.py) ... done
  Created wheel for bayesian-optimization: filename=bayesian_optimization-1.0.1-cp36-none-any.whl size=10032 sha256=d476f4cd72ce4ba849bb66887908a3f60d92da283cb0b8d59444dd266ba3fcc4
  Stored in directory: /root/.cache/pip/wheels/1d/0d/3b/6b9d4477a34b3905f246ff4e7acfd6aafd4cc9b77d473629b77
Successfully built bayesian-optimization
Installing collected packages: bayesian-optimization
Successfully installed bayesian-optimization-1.0.1
```

```

In [3]: 1 #https://www.kaggle.com/returnofspudnik/bayes-opt-again/code
2 #https://www.dlology.com/blog/how-to-do-hyperparameter-search-with-bays
3 import lightgbm as lgb
4 from bayes_opt import BayesianOptimization
5 train_df = df[df['TARGET'].notnull()]
6 test_df = df[df['TARGET'].isnull()]
7 y = train_df['TARGET']
8
9 data = lgb.Dataset(data=train_df, label=y)
10
11 def parameters(num_iterations, num_leaves, feature_fraction, max_depth,
12               params = {'application': 'binary', 'early_stopping_round': 100, 'metric':
13               params['num_iterations'] = int(round(num_iterations))
14               params["num_leaves"] = int(round(num_leaves))
15               params["learning_rate"] = learning_rate
16               params['feature_fraction'] = max(min(feature_fraction, 1), 0)
17               params['max_depth'] = int(round(max_depth))
18               params['min_split_gain'] = min_split_gain
19               params['min_child_weight'] = min_child_weight
20
21               cv_result = lgb.cv(params, data, nfold=5, seed=6, stratified=True, met
22               return max(cv_result['auc-mean'])
23
24 optimizer = BayesianOptimization(parameters, {'num_iterations': (9800, 11
25                                           'num_leaves': (16, 48),
26                                           'feature_fraction': (0.1, 0.5
27                                           'max_depth': (6, 12),
28                                           'min_split_gain': (0.001, 0.1
29                                           'min_child_weight': (20, 60)
30                                           'learning_rate': (0.01, 0.07)
31                                           })
32
33 optimizer.maximize(init_points=30, n_iter=30)
34
35 print(optimizer.max)
36

```

	iter	target	featur...	learni...	max_depth	min_ch...
	min_sp...	num_it...	num_le...			
	1	1.0	0.1895	0.03089	9.253	28.32
0.06894		1.078e+0	28.33			
	2	1.0	0.4026	0.05044	6.33	39.45
0.09028		1.051e+0	26.69			
	3	1.0	0.2039	0.03873	11.05	22.67
0.005001		1.098e+0	26.39			
	4	1.0	0.1684	0.03317	6.978	49.63
0.02587		9.879e+0	27.18			
	5	1.0	0.2845	0.03154	6.63	51.96
0.0845		1.017e+0	33.57			
	6	1.0	0.2323	0.04324	9.304	24.97
0.03821		1.06e+04	28.55			
	7	1.0	0.1349	0.0657	10.06	45.71
0.09596		1.014e+0	38.51			
	8	1.0	0.4099	0.01673	7.422	31.73
0.01013		1.081e+0	23.89			
	9	1.0	0.1966	0.04361	10.71	59.45
0.07365		1.069e+0	33.45			
	10	1.0	0.3293	0.03728	6.678	33.34
0.04531		9.831e+0	17.92			
	11	1.0	0.3723	0.02842	10.25	31.12
0.08445		1.052e+0	19.29			
	12	1.0	0.1611	0.02491	9.896	34.35
0.05449		1.083e+0	27.83			
	13	1.0	0.1106	0.06462	10.03	32.55
0.02837		1.097e+0	36.68			
	14	1.0	0.2265	0.04103	8.377	46.33


```
In [0]: 1 #df = reduce_mem_usage(df)
        2 train_df = df[df['TARGET'].notnull()]
        3 test_df = df[df['TARGET'].isnull()]
```

5 k-cross validation

```
In [0]: 1 f=KFold(n_splits=5,shuffle=True,random_state=101) #K fold cross validation
        2
        3 imp = pd.DataFrame()
        4 y=train_df['TARGET']
        5
```

```
In [0]: 1 feats = [f for f in train_df.columns if f not in ['TARGET', 'SK_ID_CURR']]
        2 imp['feat']=feats
```

```
In [0]: 1 #hold outputs of training, testing and cv data
        2 train_predict = np.zeros(train_df.shape[0])
        3 cv_predict = np.zeros(train_df.shape[0])
        4 test_predict = np.zeros(test_df.shape[0])
        5
```

Training on LGBM

Advantages:

- Fast training efficiency because of histogram binning.
- Better accuracy because of leaf wise growth.
- Low memory usage because it replaces continuous values to discrete bins which result in lower memory usage

```

In [26]: 1
2 import pickle
3 a=0
4 for i,(train, cv) in enumerate(f.split(train_df[feats],y)):
5     X_train, Y_train = train_df[feats].iloc[train], y.iloc[train]
6     X_valid, Y_valid = train_df[feats].iloc[cv], y.iloc[cv]
7
8     lgb = LGBMClassifier(
9         nthread=4,
10        n_estimators=10780,
11        learning_rate=0.03,
12        num_leaves=28,
13        feature_fraction=0.1894, #lgbm will select 30 % of the featu
14        max_depth=9,
15        min_split_gain=0.06,
16        min_child_weight=28,
17        silent=1,
18        verbose=-1, )
19    lgb.fit(X_train, Y_train, eval_set=[(X_train, Y_train), (X_valid, Y_v
20        eval_metric= 'auc', verbose= 200, early_stopping_rounds= 20
21    a=a+1
22    train_predict[train] = lgb.predict_proba(X_train, num_iteration=lgb.b
23    cv_predict[cv]=lgb.predict_proba(X_valid, num_iteration=lgb.best_iter
24    test_predict += lgb.predict_proba(test_df[feats], num_iteration=lgb.b
25    with open('model'+ str(a) +'.pickle', 'wb') as handle:
26        pickle.dump(lgb, handle)
27
28    #train_fpr, train_tpr, tr_thresholds = roc_curve(Y_train, train_predi
29    #test_fpr, test_tpr, te_thresholds = roc_curve(Y_valid, cv_predict)
30
31    imp["imp"] = lgb.feature_importances_
32    imp["fold"] = i + 1
33    print('Fold ',i + 1,' TRAIN AUC : ',roc_auc_score(Y_train, train_prec
34    print('Fold ',i + 1,' CV AUC : ',roc_auc_score(Y_valid, cv_predict[cv
35
36
37    del lgb, X_train, Y_train, X_valid, Y_valid
38    gc.collect()
39
40    print('Full TRAIN AUC score ',roc_auc_score(y, train_predict))
41    print('Full CV AUC score ',roc_auc_score(y, cv_predict))
42
43    train_fpr, train_tpr, tr_thresholds = roc_curve(y, train_predict)
44    test_fpr, test_tpr, te_thresholds = roc_curve(y, cv_predict)
45
46    test_df['TARGET'] = test_predict
47    test_df[['SK_ID_CURR', 'TARGET']].to_csv('gbvh.csv', index= False)
48
49    #plot auc curve
50    plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, t
51    plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_
52    plt.legend()
53    plt.xlabel("K: hyperparameter")
54    plt.ylabel("AUC")
55    plt.title("ERROR PLOTS")
56    plt.grid()
57    plt.show()
58
59    #plot imp features
60    features = imp["feat"]
61    importances = imp["imp"]
62    indices = (np.argsort(importances))[-25:]
63    plt.figure(figsize=(10,12))
64    plt.title('Feature Importances')
65    plt.barh(range(len(indices)), importances[indices], color='r', align='c
66    plt.yticks(range(len(indices)), [features[i] for i in indices])
67    plt.xlabel('Relative Importance')
68    plt.show()

```

10 K-cross validation

```
In [0]: 1 #df = reduce_mem_usage(df)
2 train_df = df[df['TARGET'].notnull()]
3 test_df = df[df['TARGET'].isnull()]

In [0]: 1 f=KFold(n_splits=10,shuffle=True,random_state=101) #K fold cross validate
2
3 imp = pd.DataFrame()
4 y=train_df['TARGET']
5

In [0]: 1 feats = [f for f in train_df.columns if f not in ['TARGET','SK_ID_CURR']]
2 imp['feat']=feats

In [0]: 1 #hold outputs of training, testing and cv data
2 train_predict = np.zeros(train_df.shape[0])
3 cv_predict = np.zeros(train_df.shape[0])
4 test_predict = np.zeros(test_df.shape[0])
```

```

In [10]: 1
2 import pickle
3 a=0
4 for i,(train, cv) in enumerate(f.split(train_df[feats],y)):
5     X_train, Y_train = train_df[feats].iloc[train], y.iloc[train]
6     X_valid, Y_valid = train_df[feats].iloc[cv], y.iloc[cv]
7
8     lgb = LGBMClassifier(
9         nthread=4,
10        n_estimators=10780,
11        learning_rate=0.03,
12        num_leaves=28,
13        feature_fraction=0.1894, #lgbm will select 30 % of the featu
14        max_depth=9,
15        min_split_gain=0.06,
16        min_child_weight=28,
17        silent=1,
18        verbose=-1, )
19    lgb.fit(X_train, Y_train, eval_set=[(X_train, Y_train), (X_valid, Y_v
20        eval_metric= 'auc', verbose= 200, early_stopping_rounds= 20
21    a=a+1
22    train_predict[train] = lgb.predict_proba(X_train, num_iteration=lgb.b
23    cv_predict[cv]=lgb.predict_proba(X_valid, num_iteration=lgb.best_iter
24    test_predict += lgb.predict_proba(test_df[feats], num_iteration=lgb.b
25    with open('model'+ str(a) +'.pickle', 'wb') as handle:
26        pickle.dump(lgb, handle)
27
28
29    imp["imp"] = lgb.feature_importances_
30    imp["fold"] = i + 1
31    print('Fold ',i + 1,' TRAIN AUC : ',roc_auc_score(Y_train, train_pred
32    print('Fold ',i + 1,' CV AUC : ',roc_auc_score(Y_valid, cv_predict[cv
33
34
35    del lgb, X_train, Y_train, X_valid, Y_valid
36    gc.collect()
37
38    print('Full TRAIN AUC score ',roc_auc_score(y, train_predict))
39    print('Full CV AUC score ',roc_auc_score(y, cv_predict))
40
41    train_fpr, train_tpr, tr_thresholds = roc_curve(y, train_predict)
42    test_fpr, test_tpr, te_thresholds = roc_curve(y, cv_predict)
43
44    test_df['TARGET'] = test_predict
45    test_df[['SK_ID_CURR', 'TARGET']].to_csv('gbvh.csv', index= False)
46
47    #plot auc curve
48    plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, t
49    plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_
50    plt.legend()
51    plt.xlabel("K: hyperparameter")
52    plt.ylabel("AUC")
53    plt.title("ERROR PLOTS")
54    plt.grid()
55    plt.show()
56
57    #plot imp features
58    features = imp["feat"]
59    importances = imp["imp"]
60    indices = (np.argsort(importances))[-25:]
61    plt.figure(figsize=(10,12))
62    plt.title('Feature Importances')
63    plt.barh(range(len(indices)), importances[indices], color='r', align='c
64    plt.yticks(range(len(indices)), [features[i] for i in indices])
65    plt.xlabel('Relative Importance')
66    plt.show()

```

Training until validation scores don't improve for 200 rounds.

120001 training's auc: 0.805515 training's binary logloss: 0.23332

##Conclusion

```
In [0]: 1 from prettytable import PrettyTable
2
3 x = PrettyTable()
4
5 x.field_names = ["Model", "Metrics", 'test score']
6
7 x.add_row(['Logistic regression', 'AUC', 0.7480])
8 x.add_row(['XGBoost', 'AUC', 0.7710])
9 x.add_row(['Random forest', 'AUC', 0.7417])
10 x.add_row(['lightGBM', 'AUC', 0.7948,])
11 print(x)
```

Observation:

- We can see that lightgbm is the clear winner here with 79.48% AUC

Steps for model

- Download the home credit default risk dataset.
- Perform EDA and data analysis for all the files.
- Extract important features.
- Combine all the file and Prepare data for models.
- Use different models like logistic regression, xgboost, random forest, lgbm.
- Perform hyperparameter tuning to get the best parameter.
- Train models on best hyperparameters.
- Check AUC for various model and give the best model.
- LightGBM performs best among all the models.

Featurization

- First i used simple domain specific features of application data.
- Then i used aggregation(mean, median, var) on columns for all table except application table.
- I had nearly 750 features and i selected top 300 features.
- Training xgboost on 300 features gave me an AUC of 77%.
- Then i used time related features and bayesian optimization for parameter tuning and lightgbm for training.
- I had 2337 features now and i used bayesian optimization on them and got the best hyperparameters.
- I trained lightgbm on best parameters on all the features and got an AUC of 79.79% on test data.

Refrences

- Bayesian optimization
<https://www.kaggle.com/returnofspatnik/bayes-opt-again/code> (<https://www.kaggle.com/returnofspatnik/bayes-opt-again/code>)
<https://www.dlology.com/blog/how-to-do-hyperparameter-search-with-baysian-optimization-for-keras-model/> (<https://www.dlology.com/blog/how-to-do-hyperparameter-search-with-baysian-optimization-for-keras-model/>)
- Time related features
<https://www.kaggle.com/c/home-credit-default-risk/discussion/64593> (<https://www.kaggle.com/c/home-credit-default-risk/discussion/64593>)
- Domain specific features
<https://www.kaggle.com/aantonova/797-lgbm-and-bayesian-optimization> (<https://www.kaggle.com/aantonova/797-lgbm-and-bayesian-optimization>)
- Aggregation featurization
<https://www.kaggle.com/jsaguiar/lightgbm-7th-place-solution> (<https://www.kaggle.com/jsaguiar/lightgbm-7th-place-solution>)
<https://www.linkedin.com/pulse/winning-9th-place-kaggles-biggest-competition-yet-home-levinson/>
<https://www.linkedin.com/pulse/winning-9th-place-kaggles-biggest-competition-yet-home-levinson/>
- Reduce dataframe size
<https://www.kaggle.com/gemartin/load-data-reduce-memory-usage> (<https://www.kaggle.com/gemartin/load-data-reduce-memory-usage>)

AUC on the test data

Featured Prediction Competition

Home Credit Default Risk

Can you predict how capable each applicant is of repaying a loan?

\$70,000 Prize Money

Home Credit Group · 7190 teams · a year ago

Overview Data Notebooks Discussion Leaderboard Rules Team **My Submissions** Late Submission

Your most recent submission

Name	Submitted	Wait time	Execution time	Score
test1_model.csv	just now	0 seconds	0 seconds	0.79790

Complete

[Jump to your position on the leaderboard](#)

In []: 1