

In [1]: #importing different in-built python libraries  
import numpy as np #used for mathematical programming  
import pandas as pd #data restructuring  
import matplotlib.pyplot as plt #data visualization  
import seaborn as sns

In [2]: application\_record = pd.read\_csv('application\_record.csv') #storing the given data in var  
tables  
credit\_record = pd.read\_csv('credit\_record.csv')

In [3]: application\_record.head() #displaying data

	ID	CODE	GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	NAME_INCOME_T
0	5008804	M	Y	Y	Y	0	427500.0	Ver
1	5008805	M	Y	Y	Y	0	427500.0	Ver
2	5008806	M	Y	Y	Y	0	112500.0	Ver
3	5008808	F	N			0	270000.0	Commercial asso
4	5008809	F	N	Y	Y	0	270000.0	Commercial asso

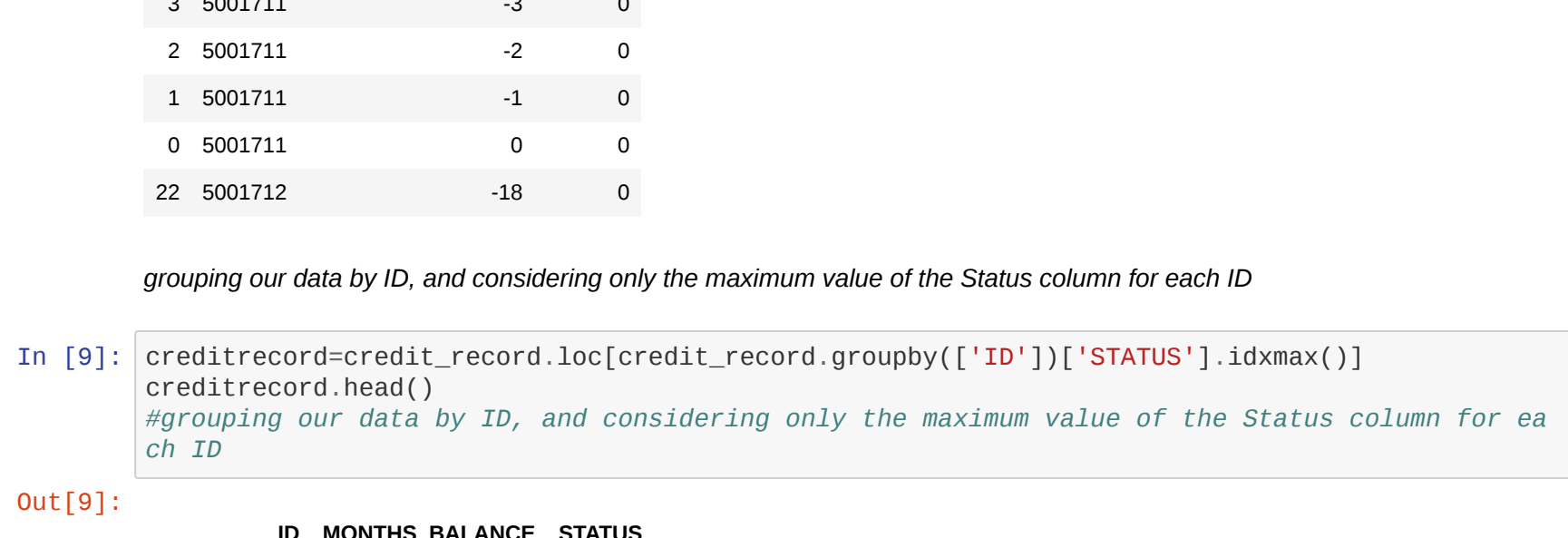
In [4]: application\_record.shape #this displays the (number of rows in our dataset, number of colum  
n in dataset)

Out[4]: (438557, 18)

In [5]: application\_record['ID'].nunique() #the number of unique IDs

Out[5]: 438510

In [6]: credit\_record.head()  
# 0: 1-29 days past due  
#1: 30-59 days past due  
#2: 60-89 days overdue  
#3: 90-119 days overdue  
#4: 120-149 days overdue  
#5: Overdue or bad debts, write-offs for more than 150 days  
#C: paid off that month X: No loan for the month



In [7]: credit\_record[credit\_record.sort\_values(['ID', 'MONTHS\_BALANCE'], ascending=True)  
credit\_record['STATUS'].replace({'C': 0, 'X': 0}, inplace=True)  
credit\_record['STATUS'] = credit\_record['STATUS'].astype('int')  
credit\_record['STATUS'] = credit\_record['STATUS'].apply(lambda x:1 if x >= 1 else 0)

```
In [10]: creditrecord.shape # (number of rows, number of columns)
Out[10]:
(45985, 3)

In [11]: application_record.info()
# there are many null values in occupation_type so we drop it

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 438557 entries, 0 to 438556
Data columns (total 10 columns):
#   #  column  non-null  dtype
#   #  -----  -
0   #  CODE     438557  non-null  int64
1   #  GENDER   438557  non-null  object
2   #  FLAG_OWN_CAR  438557  non-null  object
3   #  FLAG_OWN_REALTY  438557  non-null  object
4   #  CNT_CHILDREN  438557  non-null  int64
5   #  AMT_INCOME_TOTAL  438557  non-null  float64
6   #  NAME_INCOME_TYPE  438557  non-null  object
7   #  NAME_EDUCATION_TYPE  438557  non-null  object
8   #  NAME_FAMILY_STATUS  438557  non-null  object
```

grouping our data by ID, and considering only the maximum value of the Status column for each ID

In [9]: credit\_record=credit\_record.loc[credit\_record.groupby(['ID'])['STATUS'].idxmax()]  
credit\_record.head()

Out[9]:

```

application_record=application_record.drop('OCCUPATION_TYPE',axis=1) #drop because of NA V
# from sklearn.preprocessing import LabelEncoder
#label_encoder used to convert categorical data to numerical data
#label=LabelEncoder()
application_record.loc[:,['CODE_GENER', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY', 'NAME_INCOME_TYPE',
'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE']] = application_record[['CODE_GENER', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY', 'NAME_INCOME_TYPE',
'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE']].apply(lambda x: x.value)
application_record.DAYS_BIRTH=application_record.DAYS_BIRTH*-1

```

In [10]: credit\_record.shape # (number of rows, number of columns)

Out[10]: (45985, 3)

In [11]: application\_record.info()  
#There are many null values in occupation\_type so we drop it

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 438557 entries, 0 to 438556  
Data columns (total 18 columns):  
ID                438557 non-null int64  
CODE              438557 non-null object  
FLAG_OWN_CAR      438557 non-null object  
FLAG_OWN_REALTY   438557 non-null object  
CNT_CHILDREN      438557 non-null int64  
AMT_INCOME_TOTAL  438557 non-null float64  
NAME_INCOME_TYPE  438557 non-null object  
NAME_EDUCATION_TYPE 438557 non-null object  
NAME_FAMILY_STATUS 438557 non-null object  
NAME_HOUSING_TYPE  438557 non-null object  
DAYS_BIRTH        438557 non-null int64  
DAYS_EMPLOYED     438557 non-null int64  
FLAG_MOBILE       438557 non-null int64  
FLAG_WORK_PHONE   438557 non-null int64  
FLAG_PHONE        438557 non-null int64  
FLAG_EMAIL        438557 non-null int64  
OCCUPATION_TYPE   393354 non-null object  
CNT_FAM_MEMBERS   438557 non-null float64  
dtypes: float64(2), int64(8), object(8)  
memory usage: 60.2+ MB
```

Label encoder used to convert categorical data to numerical data

In [12]: application\_record=application\_record.drop\_duplicates('ID') #dropping duplicate IDs  
application\_record=application\_record.drop('OCCUPATION\_TYPE', axis=1) #drop because of NA V  
alues  
from sklearn.preprocessing import LabelEncoder  
#label encoder used to convert categorical data to numerical data  
le=LabelEncoder()  
application\_record.loc[:, ['CODE', 'GENDER', 'FLAG\_OWN\_CAR', 'FLAG\_OWN\_REALTY', 'NAME\_INCOME\_TYPE',  
'NAME\_FAMILY\_STATUS', 'NAME\_HOUSING\_TYPE']] =  
application\_record.loc[:, ['CODE', 'GENDER', 'FLAG\_OWN\_CAR', 'FLAG\_OWN\_REALTY', 'NAME\_INCOME\_TYPE',  
'NAME\_FAMILY\_STATUS', 'NAME\_HOUSING\_TYPE']].apply(le.fit\_transform)  
application\_record.DAYS\_BIRTH=application\_record.DAYS\_BIRTH-1  
application\_record.DAYS\_EMPLOYED=application\_record.DAYS\_EMPLOYED-1  
application\_record['NAME\_EDUCATION\_TYPE'].replace({'Lower secondary':0, 'Secondary / secondar  
y special':1, 'Incomplete higher':2, 'Higher education':3, 'Academic degree':4}, inplace=True)  
application\_record['NAME\_EDUCATION\_TYPE']=application\_record['NAME\_EDUCATION\_TYPE'].astype('int')  
application\_record=application\_record.drop(['FLAG\_MOBILE', 'axis=1'])  
application\_record['DAYS\_EMPLOYED']=application\_record['DAYS\_EMPLOYED'].apply(lambda x:0 i  
f x<0 else x)  
application\_record.head()

Out[12]:

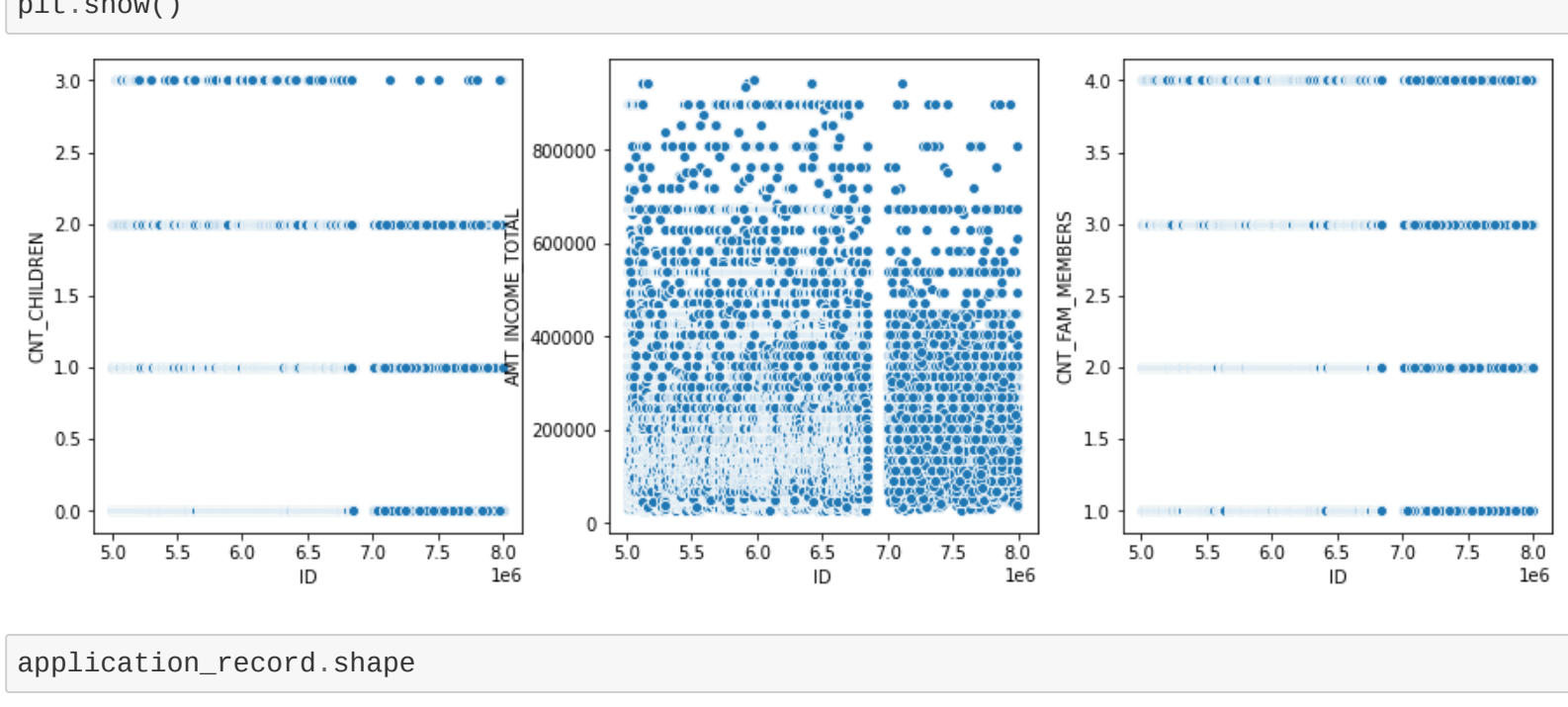
ID	CODE	GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	NAME_INCOME_T
0	5008804	1	1	1	1	0	427500.0
1	5008805	1	1	1	1	0	427500.0
2	5008806	1	1	1	1	0	112500.0
3	5008808	0	0	1	1	0	270000.0
4	5008809	0	0	1	1	0	270000.0

In [13]: application\_record.shape

Out[13]: (438510, 16)

## Outlier detection

In [14]: fig, axs = plt.subplots(3,3, figsize=(15, 5))  
sns.scatterplot(x=application\_record['ID'], y=application\_record['CNT\_CHILDREN'], ax=axs[0])  
sns.scatterplot(x=application\_record['ID'], y=application\_record['AMT\_INCOME\_TOTAL'], ax=axs[1])  
sns.scatterplot(x=application\_record['ID'], y=application\_record['CNT\_FAM\_MEMBERS'], ax=axs[2])  
plt.show()



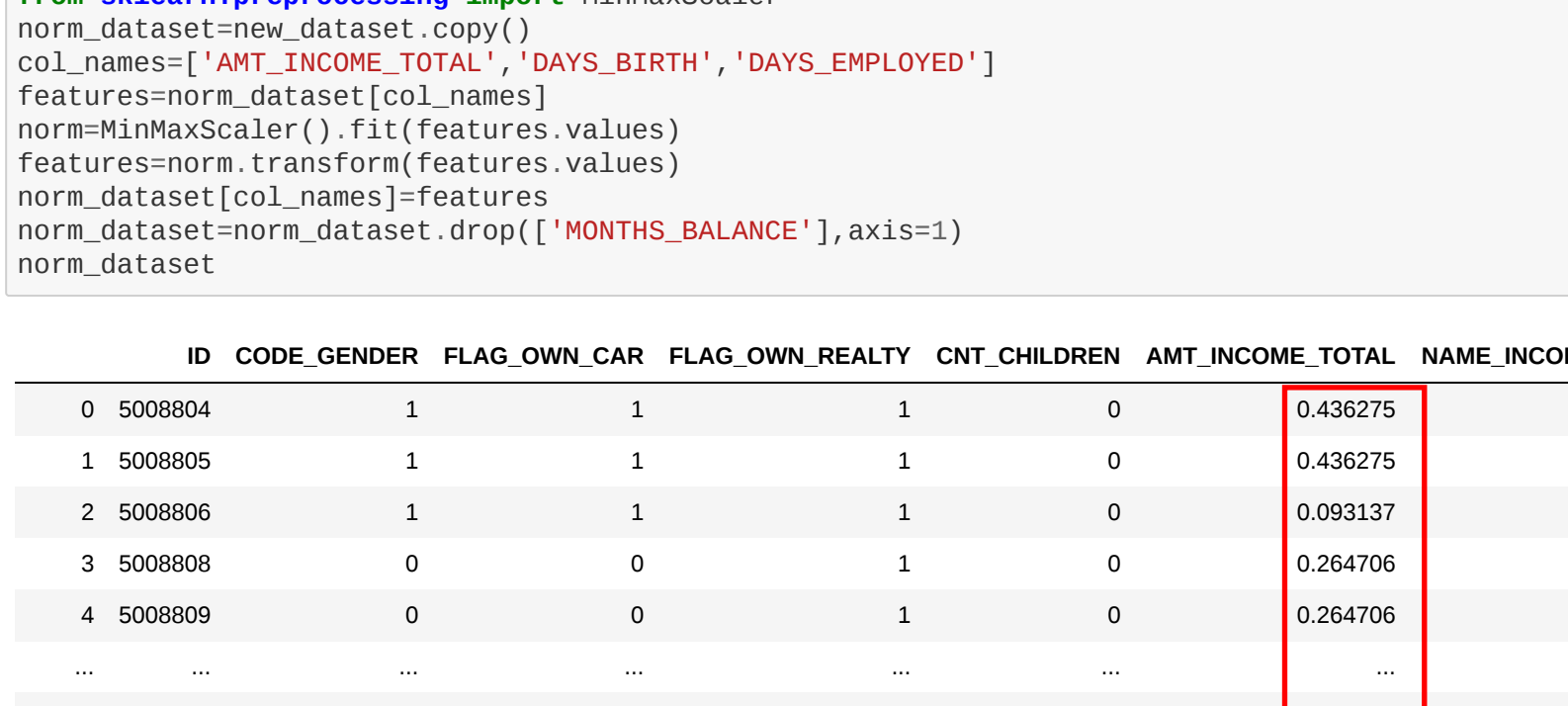
## Outlier Removal

In [15]: # We remove the top 0.1% values as we consider them to be outliers for following 3 columns  
upper\_limit=application\_record['CNT\_CHILDREN'].quantile(0.999)  
lower\_limit=application\_record['CNT\_CHILDREN'].quantile(0.000)  
application\_record = application\_record[(application\_record['CNT\_CHILDREN']>=lower\_limit) &  
(application\_record['CNT\_CHILDREN']<upper\_limit)]

In [16]: upper\_limit=application\_record['AMT\_INCOME\_TOTAL'].quantile(0.999)  
lower\_limit=application\_record['AMT\_INCOME\_TOTAL'].quantile(0.000)  
application\_record = application\_record[(application\_record['AMT\_INCOME\_TOTAL']>=lower\_limit  
& (application\_record['AMT\_INCOME\_TOTAL']<upper\_limit)]

In [17]: upper\_limit=application\_record['CNT\_FAM\_MEMBERS'].quantile(0.999)  
lower\_limit=application\_record['CNT\_FAM\_MEMBERS'].quantile(0.000)  
application\_record = application\_record[(application\_record['CNT\_FAM\_MEMBERS']>=lower\_limit  
& (application\_record['CNT\_FAM\_MEMBERS']<upper\_limit)]

In [18]: fig, axs = plt.subplots(3,3, figsize=(15, 5))  
sns.scatterplot(x=application\_record['ID'], y=application\_record['CNT\_CHILDREN'], ax=axs[0])  
sns.scatterplot(x=application\_record['ID'], y=application\_record['AMT\_INCOME\_TOTAL'], ax=axs[1])  
sns.scatterplot(x=application\_record['ID'], y=application\_record['CNT\_FAM\_MEMBERS'], ax=axs[2])  
plt.show()



In [19]: application\_record.shape

Out[19]: (432377, 16)

In [20]: new\_dataset=pd.merge(application\_record, credit\_record, on='ID')

Out[20]:

ID	CODE	GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	NAME_INCOME
0	5008804	1	1	1	1	0	427500.0
1	5008805	1	1	1	1	0	427500.0
2	5008806	1	1	1	1	0	112500.0
3	5008808	0	0	1	1	0	270000.0
4	5008809	0	0	1	1	0	270000.0

35956 rows x 18 columns

## Normalization Formula

NORMALIZATION

$$X_{new} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

In [21]: #Normalization  
from sklearn.preprocessing import MinMaxScaler  
norm\_dataset=norm\_dataset.copy()  
col\_names=['AMT\_INCOME\_TOTAL', 'DAYS\_BIRTH', 'DAYS\_EMPLOYED']  
features=norm\_dataset[col\_names]  
norm=MinMaxScaler().fit(features.values)  
features=norm.transform(features.values)  
norm\_dataset[col\_names]=features  
norm\_dataset=norm\_dataset.drop(['MONTHS\_BALANCE'], axis=1)  
norm\_dataset

Out[21]:

ID	CODE	GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	NAME_INCOME
0	5008804	1	1	1	1	0	0.436275
1	5008805	1	1	1	1	0	0.436275
2	5008806	1	1	1	1	0	0.062187
3	5008808	0	0	1	1	0	0.264706
4	5008809	0	0	1	1	0	0.264706

35956 rows x 17 columns

In [22]: #dropping ID and Status column to get our features (X)  
X\_norm=norm\_dataset.drop(['ID', 'STATUS'], axis=1)  
#Selecting the Status column to get our response variable(y)  
y\_norm\_dataset['STATUS']

In [23]: np.bincount(y\_norm\_dataset) #displays [number of 0, number of 1]

Out[23]: array([31730, 4226], dtype=int64)

High imbalance in the data

In [24]: from sklearn.model\_selection import train\_test\_split  
X\_train\_norm, X\_test\_norm, y\_train, y\_test = train\_test\_split(X\_norm, y\_norm\_dataset, ran  
dom\_state=10)

## SMOTE oversampling

In [25]: #SMOTE  
from imblearn.over\_sampling import SMOTE  
smote=SMOTE()  
X\_train\_smote, y\_train\_smote=smote.fit\_resample(X\_train\_norm, y\_train)

Using TensorFlow backend.

In [26]: np.bincount(y\_train\_smote) # y values after oversampling  
# (number of zeros, number of one)

Out[26]: array([22855, 22895], dtype=int64)

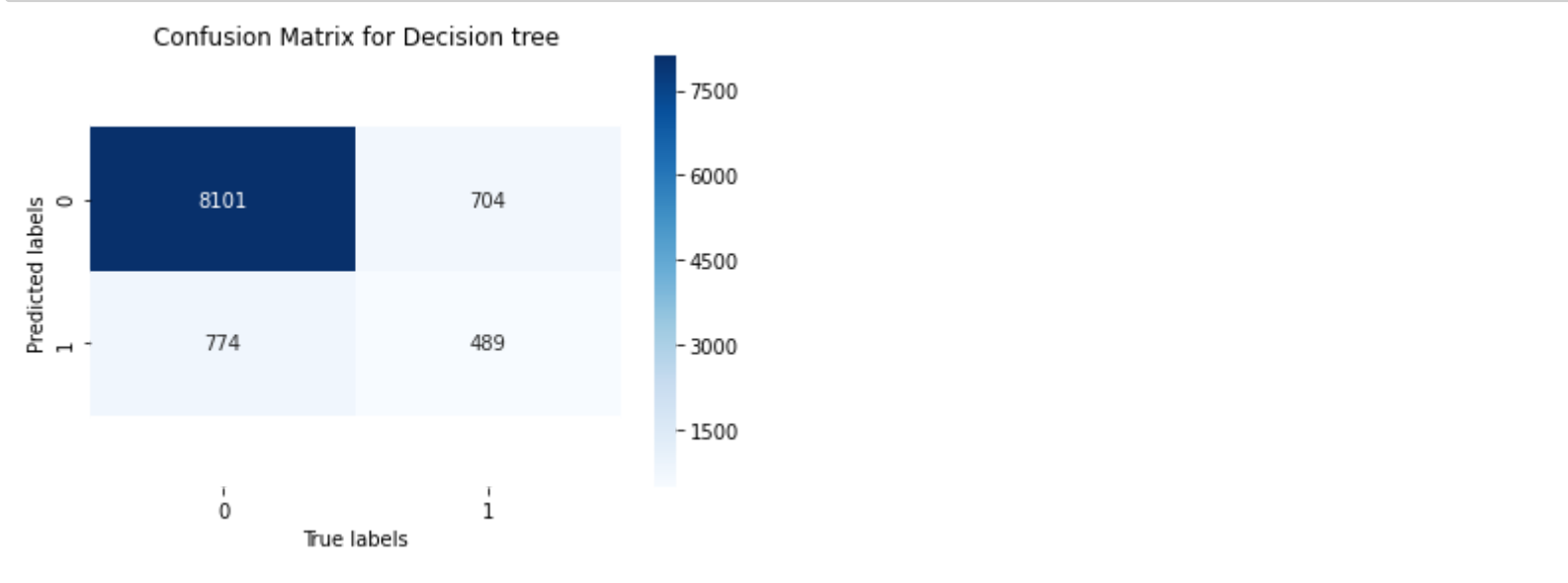
Imbalance removed from the training data using oversampling

## Logistic Regression

In [28]: from sklearn.linear\_model import LogisticRegression  
from sklearn.metrics import accuracy\_score  
from sklearn.metrics import confusion\_matrix  
clf=LogisticRegression()  
clf.fit(X\_train\_smote, y\_train\_smote)  
predictions=clf.predict(X\_test\_norm)  
print(accuracy\_score(predictions, y\_test))

C:\Users\Shravan Tawri\Anaconda3\Libsite-packages\sklearn\linear\_model\logistic.py:432: Future  
Warning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this  
warning.  
FutureWarning:

In [29]: cm=confusion\_matrix(predictions, y\_test)  
ax=plt.subplot()  
sns.heatmap(cm, annot=True, ax = ax, cmap='Blues', fmt='g')  
ax.set\_xlabel('True labels')  
ax.set\_ylabel('Predicted labels')  
ax.set\_title('Confusion Matrix for KNN')  
bottom, top = ax.get\_ylim()  
ax.set\_ylim(bottom + 0.5, top - 0.5)  
plt.show()

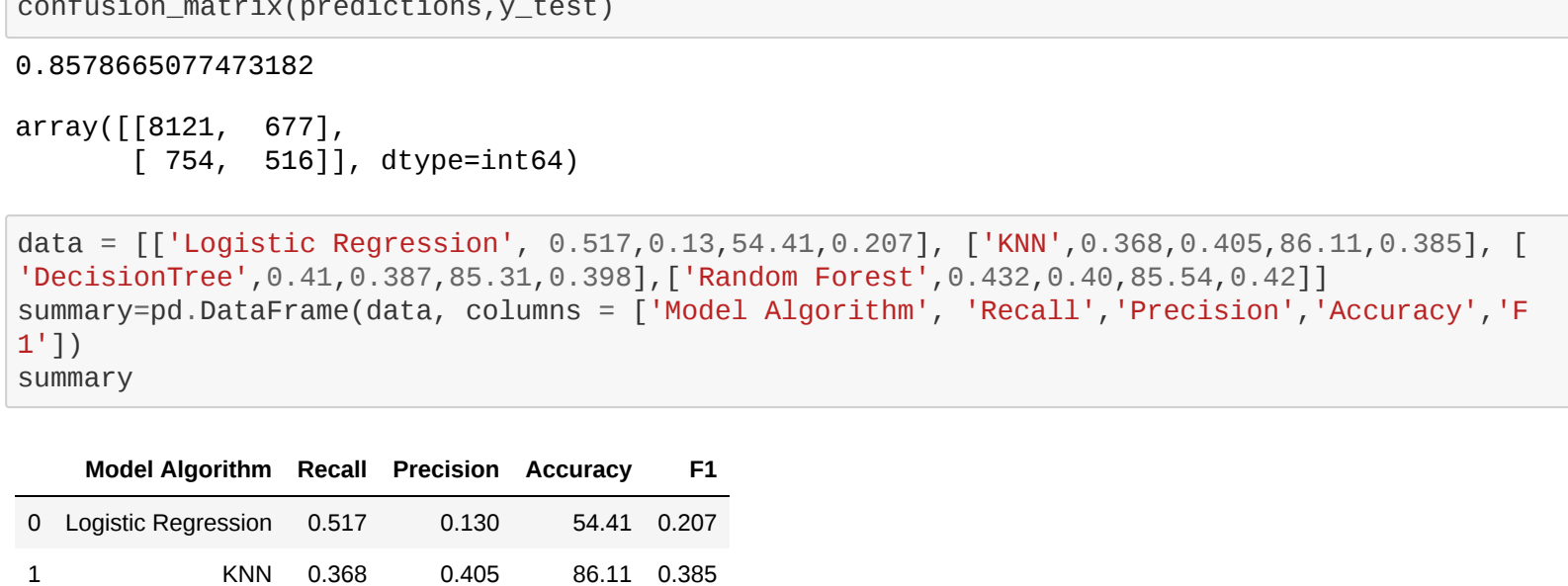


## KNN Algorithm

In [30]: from sklearn.neighbors import KNeighborsClassifier  
neigh = KNeighborsClassifier(n\_neighbors=5, weights='distance')  
neigh.fit(X\_train\_smote, y\_train\_smote)  
predictions=neigh.predict(X\_test\_norm)  
print(accuracy\_score(predictions, y\_test))

Out[30]: 0.861442193087009

In [31]: cm=confusion\_matrix(predictions, y\_test)  
ax=plt.subplot()  
sns.heatmap(cm, annot=True, ax = ax, cmap='Blues', fmt='g')  
ax.set\_xlabel('True labels')  
ax.set\_ylabel('Predicted labels')  
ax.set\_title('Confusion Matrix for KNN')  
bottom, top = ax.get\_ylim()  
ax.set\_ylim(bottom + 0.5, top - 0.5)  
plt.show()

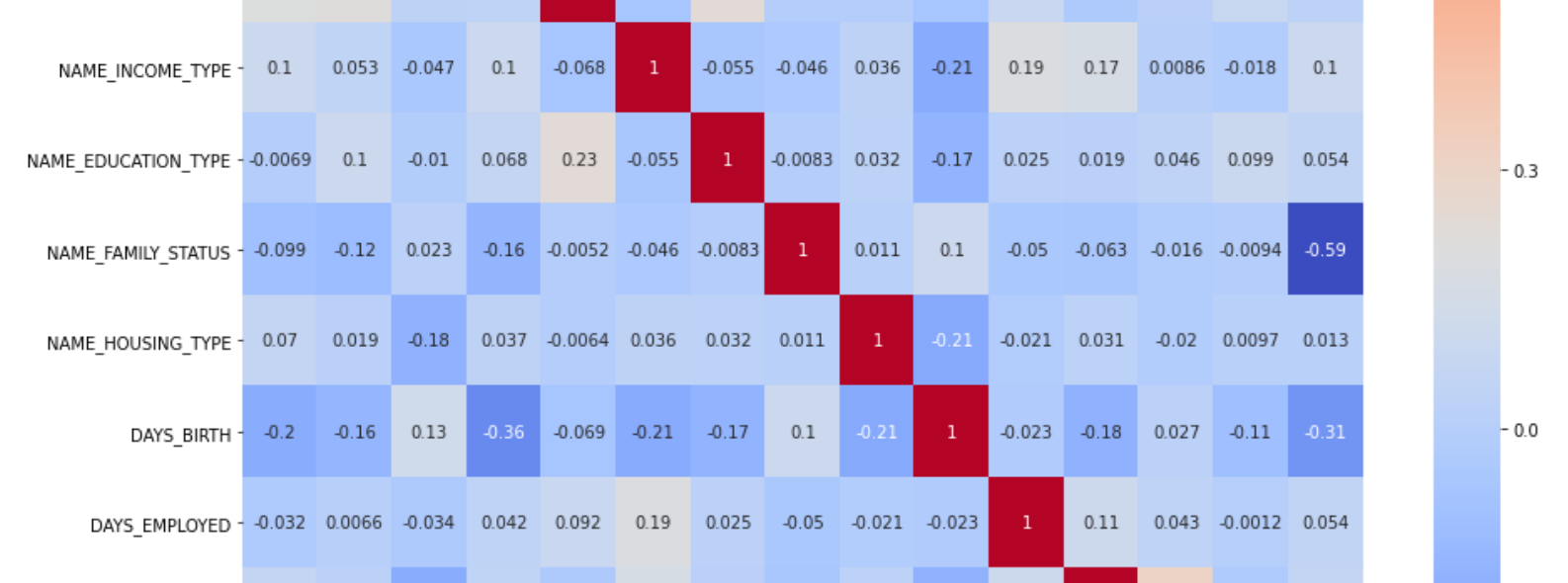


## Decision Tree Classifier

In [32]: from sklearn.tree import DecisionTreeClassifier  
clf = DecisionTreeClassifier()  
clf.fit(X\_train\_smote, y\_train\_smote)  
predictions=clf.predict(X\_test\_norm)  
print(accuracy\_score(predictions, y\_test))

Out[32]: 0.8531982518871672

In [33]: cm=confusion\_matrix(predictions, y\_test)  
ax=plt.subplot()  
sns.heatmap(cm, annot=True, ax = ax, cmap='Blues', fmt='g')  
ax.set\_xlabel('True labels')  
ax.set\_ylabel('Predicted labels')  
ax.set\_title('Confusion Matrix for Decision tree')  
bottom, top = ax.get\_ylim()  
ax.set\_ylim(bottom + 0.5, top - 0.5)  
plt.show()



## Random Forest Classifier (using GridSearch CV)

In [34]: from sklearn.ensemble import RandomForestClassifier  
clf = RandomForestClassifier()  
from sklearn.model\_selection import GridSearchCV  
parameters=[('n\_estimators', [5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100]), ('c  
riterion', ['gini', 'entropy', 'log\_loss', 'log\_likelihood', 'balanced']))  
grid\_search = GridSearchCV(estimator=clf, param\_grid=parameters, scoring='f1', cv=7)  
grid\_search.fit(X\_train\_smote, y\_train\_smote)  
grid\_search.best\_score\_

Out[34]: 0.9057837734858251

In [35]: grid\_search.best\_params\_

Out[35]: {'class\_weight': 'balanced', 'criterion': 'entropy', 'n\_estimators': 100}

In [36]: from sklearn.metrics import accuracy\_score  
clf = RandomForestClassifier()  
from sklearn.metrics import confusion\_matrix  
clf = GridSearchCV(estimator=clf, param\_grid={'n\_estimators':100, 'class\_weight':'balanced', 'criterion':'entropy'})  
clf.fit(X\_train\_smote, y\_train\_smote)  
predictions=clf.predict(X\_test\_norm)  
print(accuracy\_score(predictions, y\_test))  
confusion\_matrix(predictions, y\_test)

Out[36]: array([[ 7124, 677],  
 [ 521, 516]], dtype=int64)

In [37]: data = [['Logistic Regression', 0.537, 0.13, 54.43, 0.207], ['KNN', 0.368, 0.405, 86.11, 0.385], [  
'Decision Tree', 0.43, 0.387, 85.31, 0.398], ['Random Forest', 0.432, 0.40, 85.54, 0.42]]  
summary=pd.DataFrame(data, columns = ['Model Algorithm', 'Recall', 'Precision', 'Accuracy', 'F  
1'])  
summary

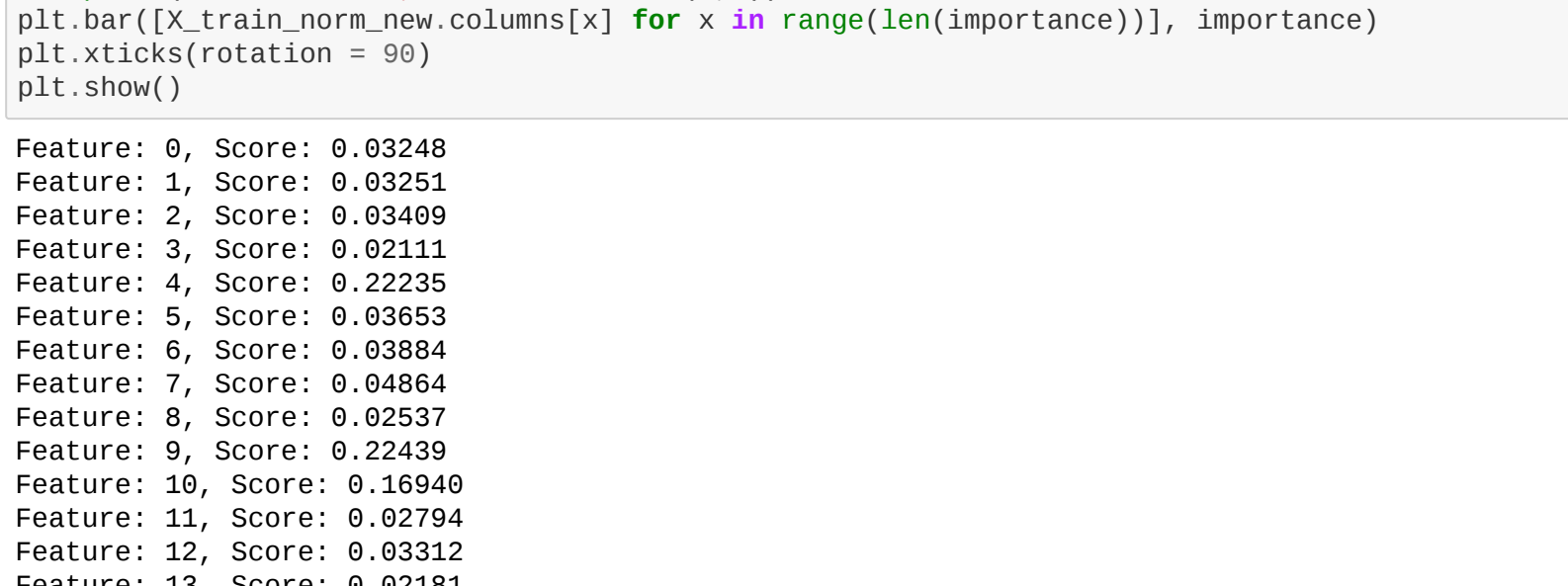
Out[37]:

Model Algorithm	Recall	Precision	Accuracy	F1
0 Logistic Regression	0.368	0.130	54.41	0.207
1 KNN	0.368	0.405	86.11	0.385
2 DecisionTree	0.410	0.387	85.31	0.398
3 Random Forest	0.432	0.400	85.54	0.420

## Correlation Matrix to check Co-linearity

In [38]: plt.figure(figsize=(15,15))  
sns.heatmap(X\_norm.corr(), annot=True, cmap='coolwarm')

Out[38]: <matplotlib.axes.\_subplots.AxesSubplot at 0x2a84bc7ec48>



removing CNT\_FAM\_MEMBERS as it is highly correlated to CNT\_CHILDREN

In [64]: X\_norm=X\_norm.drop(['ID', 'STATUS'], axis=1)  
X\_train\_norm\_new, X\_test\_norm\_new, y\_train\_new, y\_test\_new = train\_test\_split(X\_norm\_new, y,  
train\_size=0.28, random\_state=42)  
X\_train\_smote\_new, y\_train\_smote\_new=smote.fit\_resample(X\_train\_norm\_new, y\_train\_new)

In [66]: clf\_new = RandomForestClassifier(n\_estimators=100, class\_weight='balanced', criterion='entrop  
y')  
clf\_new.fit(X\_train\_smote\_new, y\_train\_smote\_new)  
predictions=clf\_new.predict(X\_test\_norm\_new)  
print(accuracy\_score(predictions, y\_test\_new))  
confusion\_matrix(predictions, y\_test\_new)

Out[66]: array([[ 7424, 673],  
 [ 521, 522]], dtype=int64)

In [68]: cm=confusion\_matrix(predictions, y\_test\_new)  
ax=plt.subplot()  
sns.heatmap(cm, annot=True, ax = ax, cmap='Blues', fmt='g')  
ax.set\_xlabel('True labels')  
ax.set\_ylabel('Predicted labels')  
ax.set\_title('Confusion Matrix')  
bottom, top = ax.get\_ylim()  
ax.set\_ylim(bottom + 0.5, top - 0.5)  
plt.show()



## Check the importance of each feature

In [71]: importance=clf\_new.feature\_importances\_

In [73]: for i, v in enumerate(importance):  
print('Feature: %4d, Score: %.5f' % (i, v))  
plt.bar([X\_train\_norm\_new.columns[x] for x in range(len(importance))], importance)  
plt.xticks(rotation = 90)

