

Credit Card Fraud Detection

Data: Number of Instances: 30000

```
In [1]: # loading library
import pandas as pd
import numpy as np

#model
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier

# Evaluation
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report

from sklearn.metrics import precision_score
from sklearn.metrics import recall_score
from sklearn.metrics import f1_score

#Visualization
import matplotlib as plt
import seaborn as sns
import matplotlib.pyplot as plt
from matplotlib import gridspec
from termcolor import colored as cl # text customization
```

Import Data

```
In [2]: # loading data
# create URL

# load dataset
dataframe = pd.read_excel("/Users/wangtingting/Documents/763-Project/default of credit card clients.xls",
                           sheet_name = 0, header = 1)

# view first two row
dataframe.head(10)
```

Out [2]:

	ID	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	...	BILL_AMT4	BILL_AMT5	BILL_AMT6	PAY_AM
0	1	20000	2	2	2	1	24	2	2	-1	-1	...	0	0	0
1	2	120000	2	2	2	2	26	-1	2	0	0	...	3272	3455	3261
2	3	90000	2	2	2	1	34	0	0	0	0	...	14331	14948	15549
3	4	50000	2	2	2	1	37	0	0	0	0	...	28314	28959	29547
4	5	50000	1	2	1	57	-1	0	0	-1	0	...	20940	19146	19131
5	6	50000	1	1	2	37	0	0	0	0	0	...	19394	19619	20024
6	7	50000	1	1	2	29	0	0	0	0	0	...	542953	483003	479944
7	8	100000	2	2	2	23	0	-1	-1	0	0	...	221	-159	567
8	9	140000	2	3	1	28	0	0	2	0	0	...	12211	11790	3719
9	10	20000	1	3	2	35	-2	-2	-2	-2	...	0	13007	13912	

10 rows x 25 columns

```
In [3]: dataframe.shape
```

Out [3]: (30000, 25)

```
In [4]: dataframe.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 25 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   ID                   30000 non-null  int64
 1   LIMIT_BAL            30000 non-null  int64
 2   SEX                  30000 non-null  int64
 3   EDUCATION            30000 non-null  int64
 4   MARRIAGE              30000 non-null  int64
 5   AGE                   30000 non-null  int64
 6   PAY_0                 30000 non-null  int64
 7   PAY_2                 30000 non-null  int64
 8   PAY_3                 30000 non-null  int64
 9   PAY_4                 30000 non-null  int64
10  PAY_5                 30000 non-null  int64
11  PAY_6                 30000 non-null  int64
12  BILL_AMT1             30000 non-null  int64
13  BILL_AMT2             30000 non-null  int64
14  BILL_AMT3             30000 non-null  int64
15  BILL_AMT4             30000 non-null  int64
16  BILL_AMT5             30000 non-null  int64
17  BILL_AMT6             30000 non-null  int64
18  PAY_AMT1              30000 non-null  int64
19  PAY_AMT2              30000 non-null  int64
20  PAY_AMT3              30000 non-null  int64
21  PAY_AMT4              30000 non-null  int64
22  PAY_AMT5              30000 non-null  int64
23  PAY_AMT6              30000 non-null  int64
24  default payment next month  30000 non-null  int64
dtypes: int64(25)
memory usage: 5.5 MB
```

```
In [5]: dataframe.columns
```

Out [5]: Index(['ID', 'LIMIT_BAL', 'SEX', 'EDUCATION', 'MARRIAGE', 'AGE', 'PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6', 'BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_AMT4', 'BILL_AMT5', 'BILL_AMT6', 'PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY_AMT6', 'default payment next month'], dtype='object')

```
In [6]: # remove useless column
data_after_drop=dataframe.drop(['ID'], axis=1)
#there are 15 columns left
data_after_drop.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 30000 entries, 0 to 29999
Data columns (total 24 columns):
 #   Column              Non-Null Count  Dtype
---  --
 0   LIMIT_BAL            30000 non-null  int64
 1   SEX                  30000 non-null  int64
 2   EDUCATION            30000 non-null  int64
 3   MARRIAGE              30000 non-null  int64
 4   AGE                   30000 non-null  int64
 5   PAY_0                 30000 non-null  int64
 6   PAY_2                 30000 non-null  int64
 7   PAY_3                 30000 non-null  int64
 8   PAY_4                 30000 non-null  int64
 9   PAY_5                 30000 non-null  int64
10  PAY_6                 30000 non-null  int64
11  BILL_AMT1             30000 non-null  int64
12  BILL_AMT2             30000 non-null  int64
13  BILL_AMT3             30000 non-null  int64
14  BILL_AMT4             30000 non-null  int64
15  BILL_AMT5             30000 non-null  int64
16  BILL_AMT6             30000 non-null  int64
17  PAY_AMT1              30000 non-null  int64
18  PAY_AMT2              30000 non-null  int64
19  PAY_AMT3              30000 non-null  int64
20  PAY_AMT4              30000 non-null  int64
21  PAY_AMT5              30000 non-null  int64
22  PAY_AMT6              30000 non-null  int64
23  default payment next month  30000 non-null  int64
dtypes: int64(24)
memory usage: 5.5 MB
```

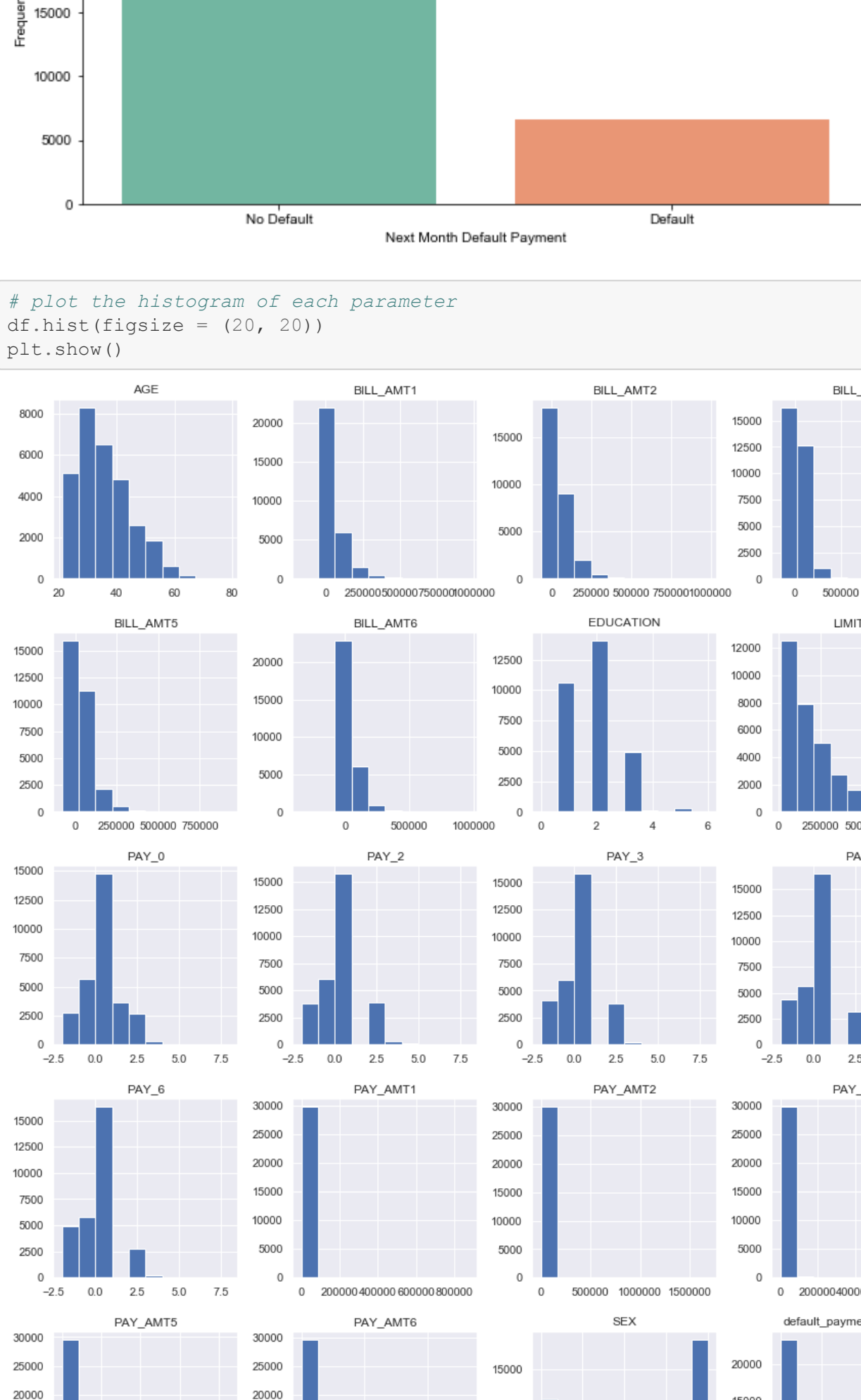
```
In [7]: # check missing values
data_after_drop[data_after_drop['default payment next month'].isnull()].count()

Out [7]:
LIMIT_BAL      0
SEX             0
EDUCATION       0
MARRIAGE        0
AGE             0
PAY_0           0
PAY_2           0
PAY_3           0
PAY_4           0
PAY_5           0
PAY_6           0
BILL_AMT1       0
BILL_AMT2       0
BILL_AMT3       0
BILL_AMT4       0
BILL_AMT5       0
BILL_AMT6       0
PAY_AMT1        0
PAY_AMT2        0
PAY_AMT3        0
PAY_AMT4        0
PAY_AMT5        0
PAY_AMT6        0
default payment next month    0
dtypes: int64
```

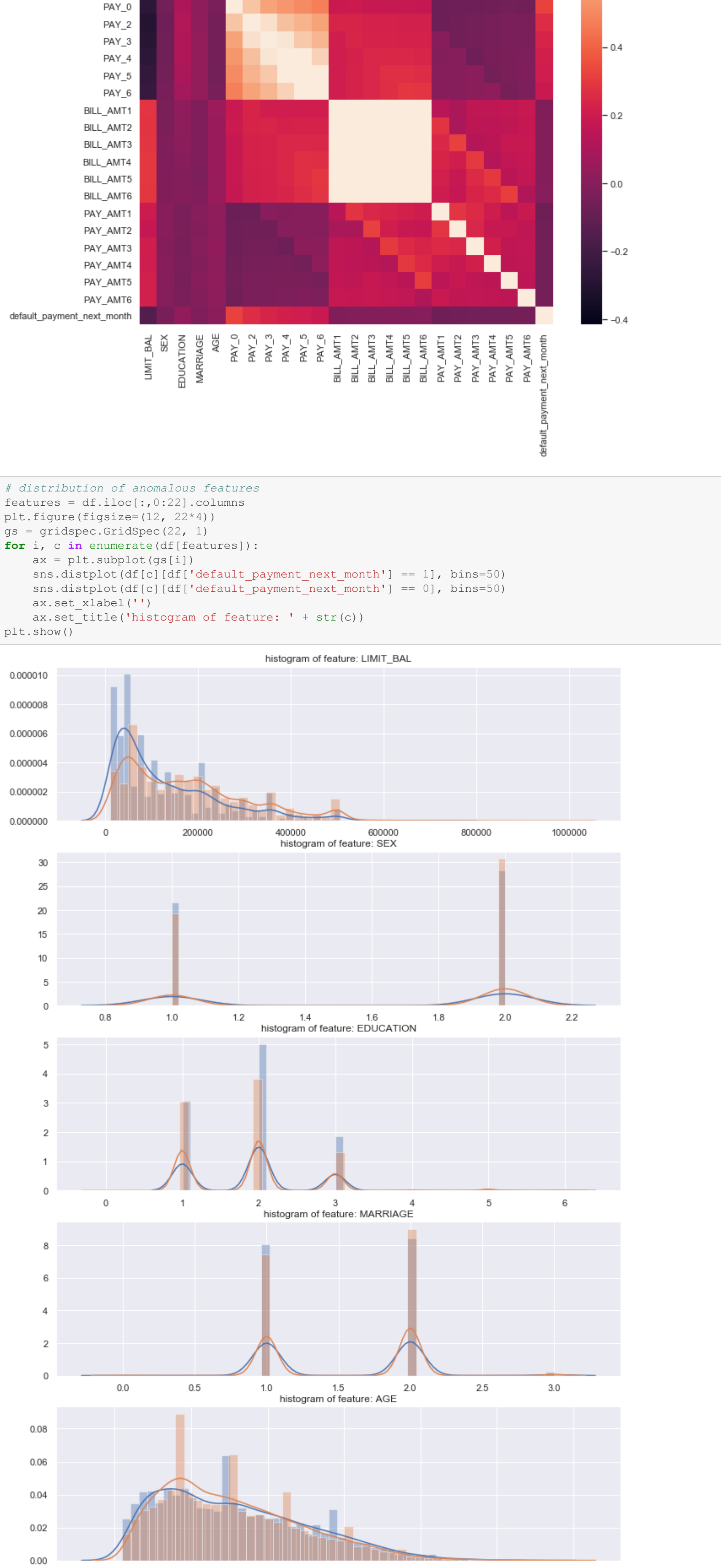
```
In [8]: # rename target column
df = data_after_drop.rename(columns={'default payment next month':'default_payment_next_month'})

#concatenate dataframe by column
df1 = df.copy()
```

```
In [9]: # check target balance or unbalance?
plt.title('Default Next Month Payment')
ax = sns.countplot(x = df.default_payment_next_month, palette="Set2")
ax.set(font_scale=1.0)
ax.set_ylim(top = 30000)
ax.set_xticklabels(['No Default', 'Default'])
ax.set_xlabel('Next Month Default Payment')
ax.set_ylabel('Frequency')
fig = plt.gcf()
fig.set_size_inches(10,5)
plt.show()
```



```
In [10]: # plot the histogram of each parameter
df.hist(figsize = (20, 20))
plt.show()
```

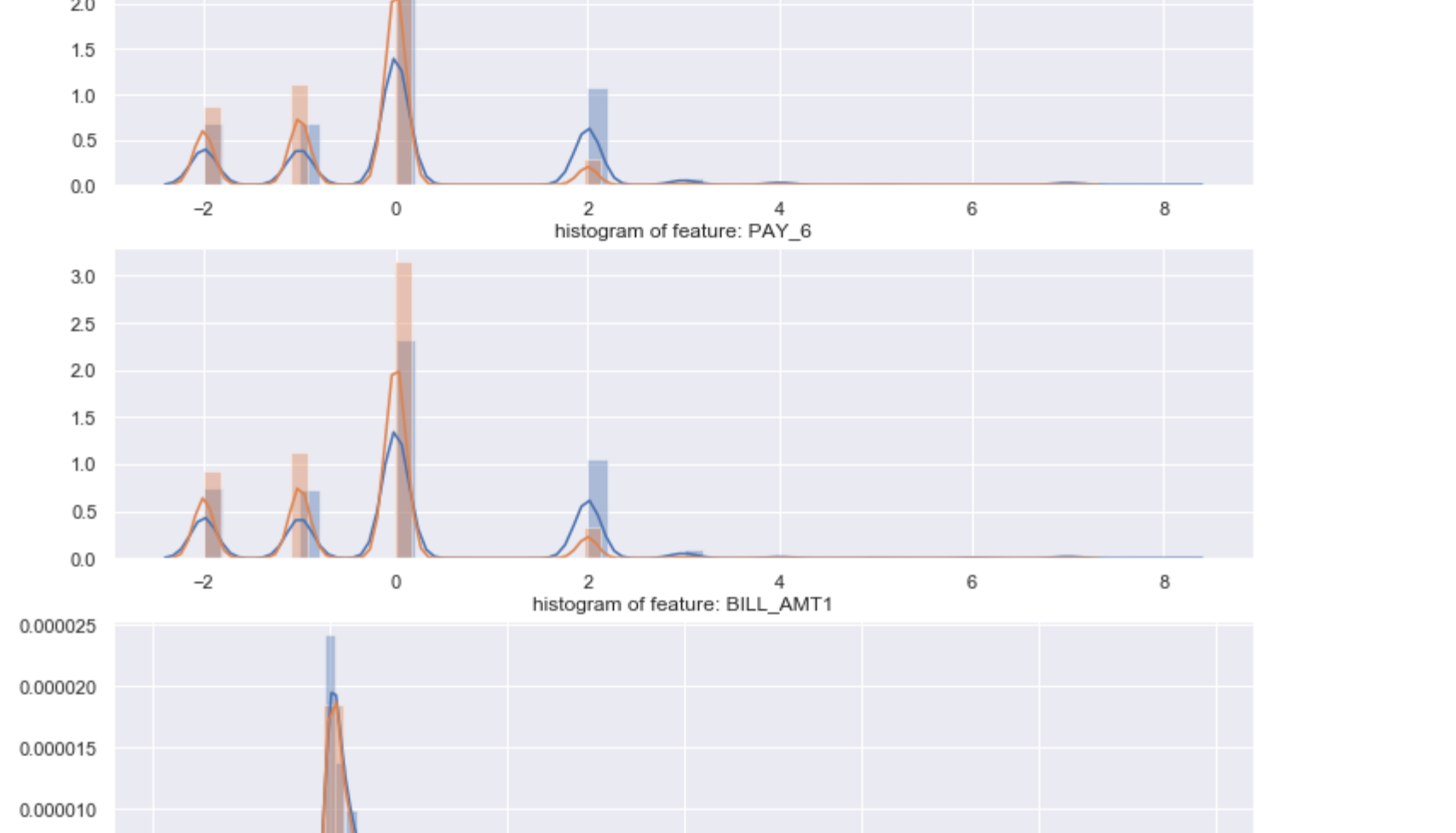


```
In [11]: # default pay
df[['PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6']].describe()
```

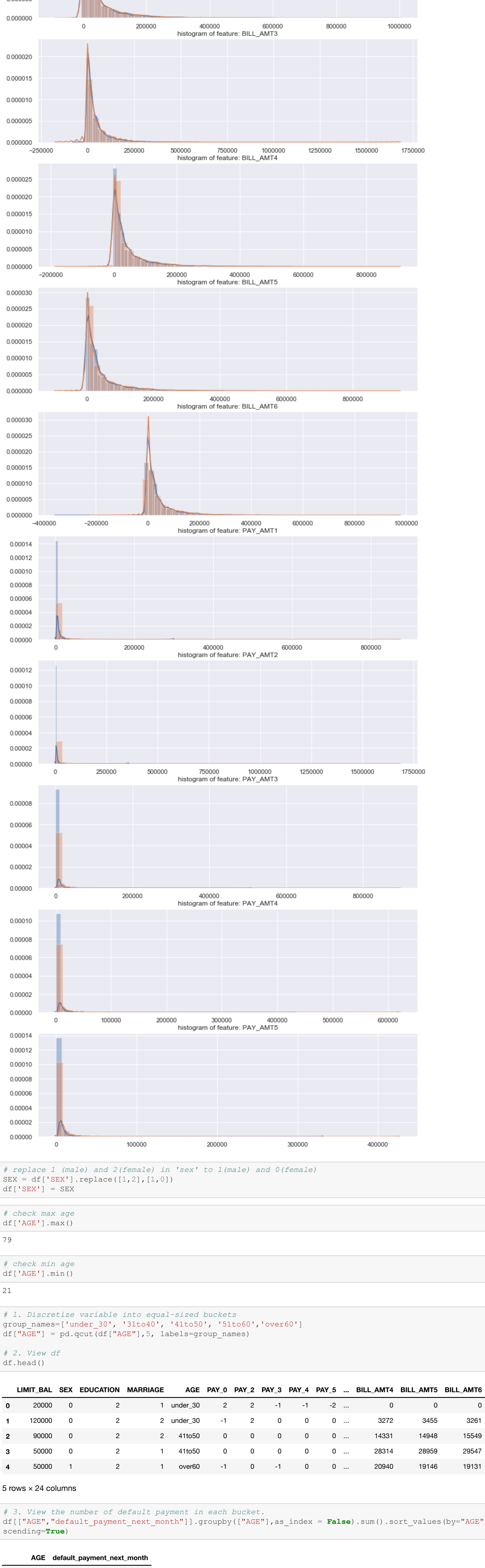
Out [11]:

	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6
count	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000
mean	-0.016700	-0.135767	-0.166300	-0.220667	-0.266200	-0.291100
std	1.123802	1.197186	1.196888	1.169139	1.133167	1.149988
min	-2.000000	-2.000000	-2.000000	-2.000000	-2.000000	-2.000000
25%	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000
50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	8.000000	8.000000	8.000000	8.000000	8.000000	8.000000

```
In [12]: # Correlation matrix
corrmat = df.corr()
fig = plt.figure(figsize = (12, 9))
sns.heatmap(corrmat, vmax = -.8, square = True)
plt.show()
```



```
In [13]: # distribution of malicious features
features = df.iloc[:,0:22].columns
plt.figure(figsize=(12, 22*4))
gs = gridspec.GridSpec(22, 1)
for i, c in enumerate(features):
    ax = plt.subplot(gs[i])
    sns.distplot(df[c][df['default_payment_next_month'] == 1], bins=50)
    sns.distplot(df[c][df['default_payment_next_month'] == 0], bins=50)
    ax.set_xlabel('')
    ax.set_title('histogram of feature: ' + str(c))
plt.show()
```



```
In [14]: # replace 1 (male) and 2(female) in 'sex' to 1(male) and 0(female)
df['sex'].replace([1,2], [0,1])
df['SEX'] = SEX
```

```
In [15]: # check max age
df['AGE'].max()
```

Out [15]: 79

```
In [16]: # check min age
df['AGE'].min()
```

Out [16]: 21

```
In [17]: # 1. Discretize variable into equal-sized buckets
group_names=['under_30', '31to40', '41to50', '51to60', 'over60']
df['AGE'] = pd.qcut(df['AGE'],5, labels=group_names)

# 2. View df
df.head()
```

Out [17]:

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	...	BILL_AMT4	BILL_AMT5	BILL_AMT6	PAY_AM
0	20000	0	2	1	under_30	2	2	-1	-1	-2	-2	...	3913	3102	...	0
1	120000	0	2	2	under_30	-1	2	0	0	0	2	...	2682	1725	...	0
2	90000	0	2	2	41to50	0	0	0	0	0	0	...	29239	14027	...	0
3	50000	0	2	1	41to50	0	0	0	0	0	0	...	46990	48233	...	0
4	50000	1	2	1	over60	-1	0	-1	0	0	0	...	8617	5670	...	0

5 rows x 24 columns

```
In [18]: # 3. View the number of default payment in each bucket.
df['AGE'].value_counts()['default_payment_next_month'].groupby(['AGE']).sum().sort_values(ascending=False)
```

Out [18]:

AGE	default_payment_next_month
under_30	1598
31to40	1102
41to50	1380
51to60	1100
over60	1456

```
In [19]: # Age into dummy variables
dum_df_age = pd.get_dummies(df['AGE'])

# merge with main df on key values
data_after_dummy_age = pd.concat([df, dum_df_age], axis=1)

# drop AGE column
data_after_dummy_age=data_after_dummy_age.drop(['AGE'], axis=1)

# drop one dummy variable 'over60'
data_after_dummy_age=data_after_dummy_age.drop(['over60'], axis=1)

data
```

Out [19]:

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	...	PAY_AMT2	PAY_AMT3	PAY_AMT4	PAY_AMT5	PAY_AMT6
0	20000	0	2	1	2	2	-1	-1	-2	-2	...	689	0	0	0	0
1	120000	0	2	2	-1	2	0	0	0	2	...	1000	1000	1000	1000	1000
2	90000	0	2	2	0	0	0	0	0	0	...	1500	1000	1000	1000	1000
3	50000	0	2	1	0	0	0	0	0	0	...	2019	1200	1100	1000	1000
4	50000	1	2	1	-1	0	-1	0	0	0	...	36681	10000	9000	10000	9000

30000 rows x 27 columns

```
In [20]: # Married into dummy variables
data_after_dummy_EDUCATION = pd.concat([data_0, dum_df], axis=1)
df1['AGE'] = 'default_payment_next_month'
data_0 = data_0.drop(['MARRIAGE'], axis=1)
data
```

Out [20]:

	LIMIT_BAL	SEX	MARRIAGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	BILL_AMT1	BILL_AMT2	...	31to40	41to50	51to60	EDUCATION
0	20000	0	1	2	2	-1	-1	-2	-2	3913	3102	...	0	0	0	0
1	120000	0	2	-1	2	0	0	0	2	2682	1725	...	0	0	0	0
2	90000	0	2	0	0	0	0	0	0	29239	14027	...	0	1	0	0
3	50000	0	1	0	0	0	0	0	0	46990	48233	...	0	1	0	0
4	50000	1	1	-1	0	-1	0	0	0	8617	5670	...	0	0	0	0

30000 rows x 33 columns

```
In [21]: # dummy variable = MARRIAGE column
data_0['MARRIAGE'] = pd.concat([data_0, dum_df], axis=1)
dum_df = pd.get_dummies(data_0['MARRIAGE'], prefix='MARRIAGE')
data
```

Out [21]:

	MARRIAGE_0	MARRIAGE_married	MARRIAGE_others	MARRIAGE_single
0	0	0	1	0
1	0	0	0	1
2	0	0	0	1
3	0	1	0	0
4	0	1	0	0

...

	MARRIAGE_0	MARRIAGE_married	MARRIAGE_others	MARRIAGE_single
29995	0	1	0	0
29996	0	0	0	1
29997	0	0	0	1
29998	0	1	0	0
29999	0	1	0	0

30000 rows x 4 columns

```
In [22]: # merge with main df on key values
data_after_dummy_MARRIAGE = pd.concat([data_0, dum_df], axis=1)
data_0 = data_0.drop(['MARRIAGE'], axis=1)
data
```

Out [22]:

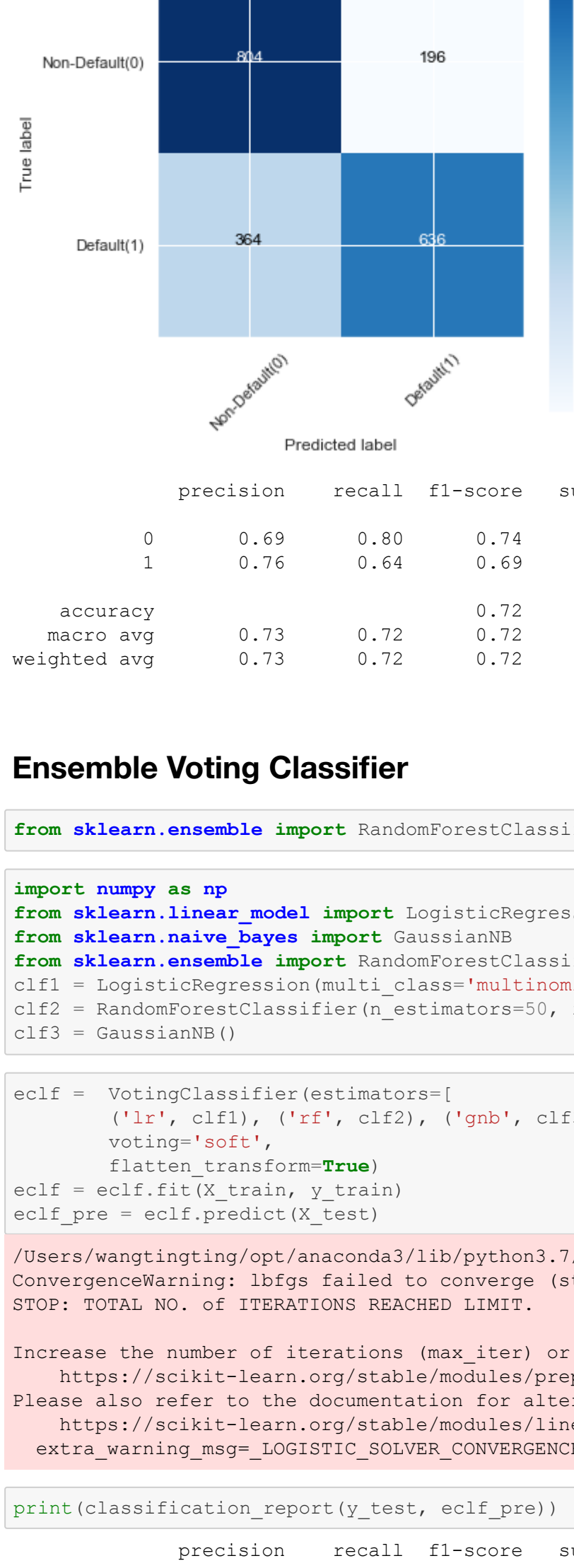
	LIMIT_BAL	SEX	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	BILL_AMT1	BILL_AMT2	...	31to40	41to50	51to60	EDUCATION
0	20000	0	2	2	-1	-1	-2	-2	3913	3102	...	0	0	0	0
1	120000	0	-1	2	0	0	0	2	2682	1725	...	0	0	0	0
2	90000	0	0	0	0	0	0	0	29239	14027	...	0	1	0	0
3	50000	0	0	0	0	0	0	0	46990	48233	...	0	1	0	0
4	50000	1	-1	0	-1	0	0	0	8617	5670	...	0	0	0	0

...

	LIMIT_BAL	SEX	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	BILL_AMT1	BILL_AMT2	...	31to40	41to50	51to60	EDUCATION
29995	220000	1	1	0	0	0	0	0	0	185948	19215	...	0	0	1
29996	150000	1	-1	-1	-1	-1	-1	0	0	1683	1828	...	0	0	1
29997	30000	1	4	3	2	-1	0	0	0	3565	3356	...	0	1	0
29998	80000	1	1	-1	0	0	0	-1	-1645	78379	...	0	0	1	0
29999	50000	1	0	0	0	0	0	0	47929	48905	...	0	0	0	0

30000 rows x 32 columns


```
[54]: # 5. Random forest tree
rf_cm_plot = plot_confusion_matrix(rf_matrix,
                                  classes = ('Not Default(0)', 'Default(1)'),
                                  normalize = False, title = 'Random Forest Tree')
plt.savefig('rf_cm_plot.png')
plt.show()
print(classification_report(y_test, rf_yhat))
```



Ensemble Voting Classifier

```
In [55]: from sklearn.ensemble import RandomForestClassifier, VotingClassifier
```

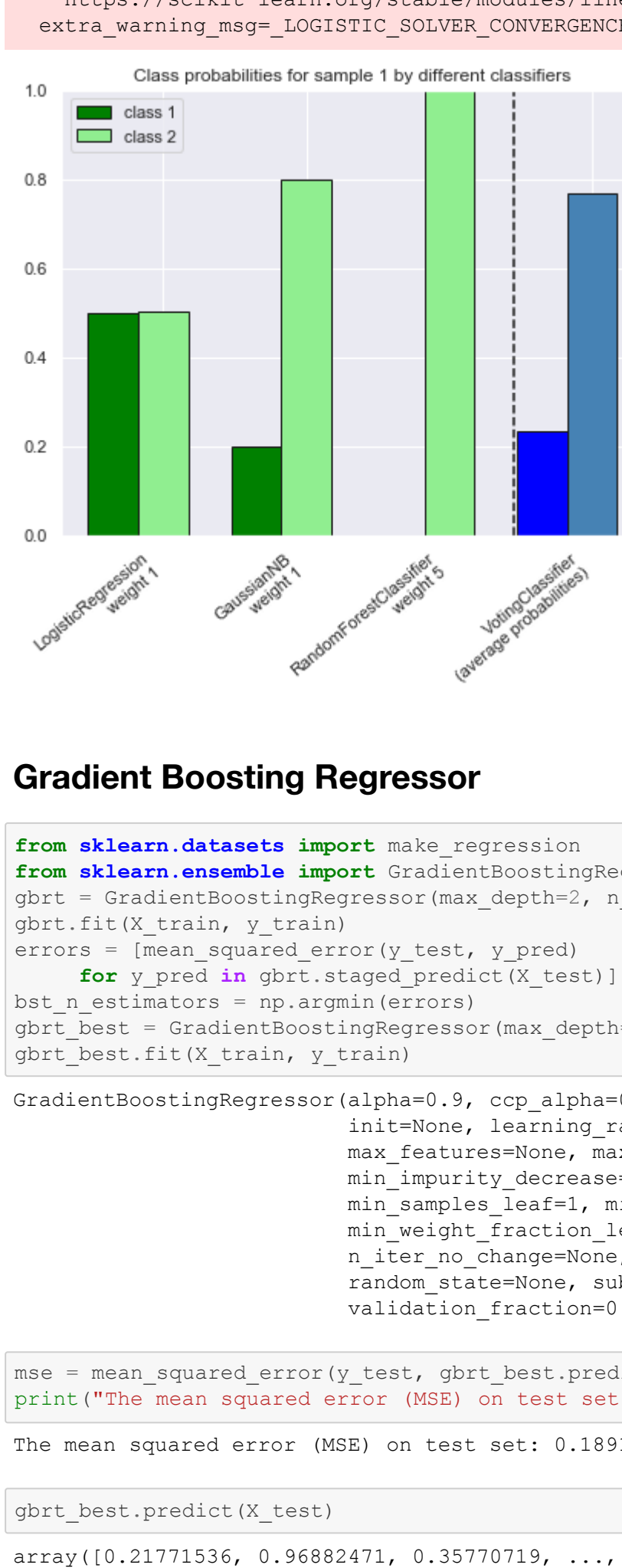
```
In [56]: import numpy as np
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import RandomForestClassifier, VotingClassifier
clf1 = LogisticRegression(multi_class='multinomial', random_state=1)
clf2 = RandomForestClassifier(n_estimators=50, random_state=1)
clf3 = GaussianNB()
```

```
In [57]: ecclf = VotingClassifier(estimators=[
    ('l1', clf1), ('rf', clf2), ('gnb', clf3)],
    voting='soft',
    flatten_transform=True)
ecclf = ecclf.fit(X_train, y_train)
ecclf_pre = ecclf.predict(X_test)
```

/Users/wangtingting/opt/anaconda3/lib/python3.7/site-packages/sklearn/linear_model/_logistic.py:940: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

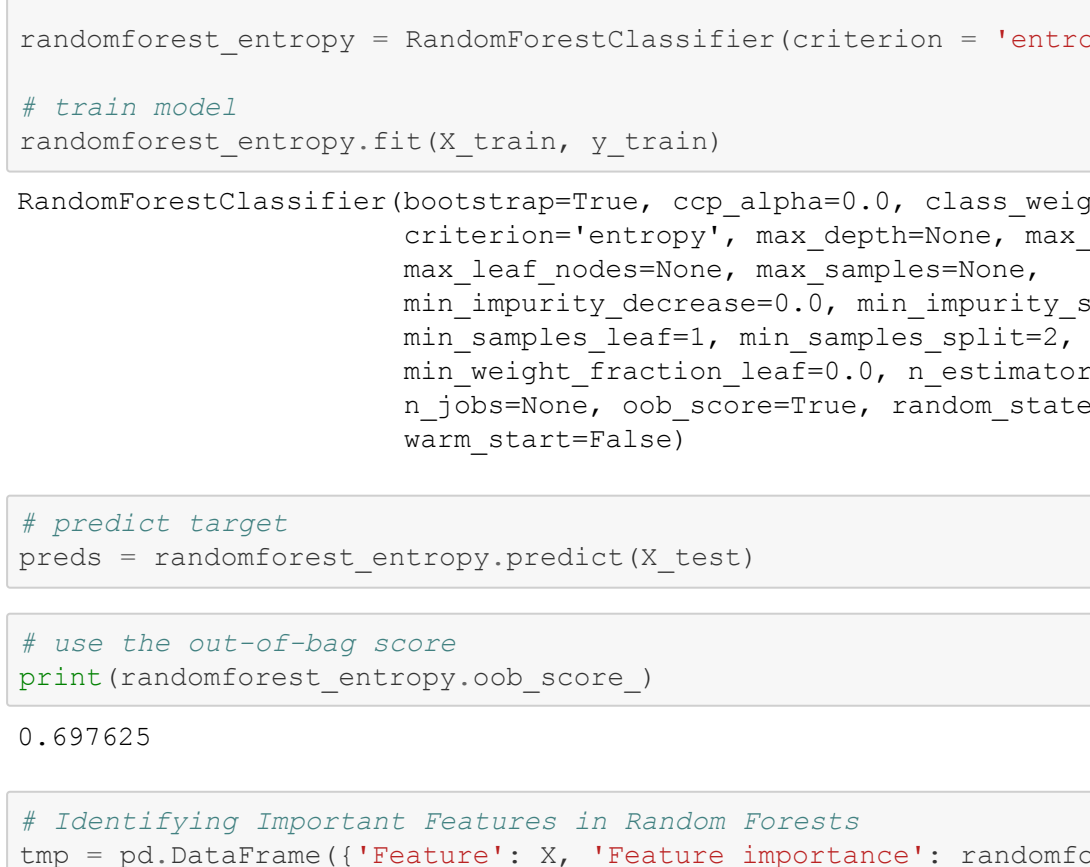
Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html>
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
extra_warning_msg="LOGISTIC_SOLVER_CONVERGENCE_MSG")

```
In [58]: print(classification_report(y_test, ecclf_pre))
```



```
In [59]: # predict class probabilities for all classifiers
probas = [c.fit(X, y).predict_proba(X) for c in (clf1, clf2, clf3, ecclf)]

# get class probabilities for the first sample in the dataset
class1_1 = [pr[0, 0] for pr in probas]
class2_1 = [pr[0, 1] for pr in probas]
```



Gradient Boosting Regressor

```
In [60]: from sklearn.datasets import make_regression
from sklearn.ensemble import GradientBoostingRegressor
```

```
gbrt = GradientBoostingRegressor(max_depth=2, n_estimators=120)
gbrt.fit(X_train, y_train)
errors = (mean_squared_error(y_test, y_pred)
          for y_pred in gbrt.staged_predict(X_test))
bat_n_estimators = np.argmax(errors)
```

```
Out[60]: GradientBoostingRegressor(alpha=0.9, ccp_alpha=0.0, criterion='friedman_mse',
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=109,
n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, tol=0.0001,
validation_fraction=0.1, verbose=0, warm_start=False)
```

```
In [61]: mse = mean_squared_error(y_test, gbrt_best.predict(X_test))
print("The mean squared error (MSE) on test set: {:.4f}".format(mse))
```

The mean squared error (MSE) on test set: 0.1891

```
In [80]: gbrt_best.predict(X_test)
```

```
Out[80]: array([0.21771536, 0.96882471, 0.35770719, ..., 0.30222342, 0.40238925,
0.46177777])
```

MODEL - Random Forest

```
In [62]: # training a random forest classifier
from sklearn.ensemble import RandomForestClassifier
```

```
In [63]: # Create random forest classifier object using entropy
```

```
randomforest_entropy = RandomForestClassifier(criterion = 'entropy', random_state=0, oob_score = True)
# train model
randomforest_entropy.fit(X_train, y_train)
```

```
Out[63]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
criterion='entropy', max_depth=None, max_features='auto',
max_leaf_nodes=None, max_samples=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=100,
n_iter_no_change=None, presort='deprecated',
random_state=None, subsample=1.0, tol=0.0001,
validation_fraction=0.1, verbose=0, warm_start=False)
```

```
In [64]: # predict target
preds = randomforest_entropy.predict(X_test)
```

```
In [65]: # use the out-of-bag score
print(randomforest_entropy.oob_score_)
```

0.697625

```
In [66]: # Identifying Important Features in Random Forests
tmp = pd.DataFrame({'Feature': X, 'Feature importance': randomforest_entropy.feature_importances_})
fig, ax1 = tmp.sort_values(by='Feature importance', ascending=False)
plt.figure(figsize = (12,8))
plt.title('Features importance', fontsize=14)
s = sns.barplot(x='Feature', y='Feature importance', data=tmp)
s.set_xticklabels(s.get_xticklabels(), rotation=90)
plt.show()
```



```
In [67]: # check CFM
from sklearn.metrics import classification_report, confusion_matrix
```

```
print(confusion_matrix(y_test, preds))
print('\n')
print(classification_report(y_test, preds))
```



```
In [68]: cm = pd.crosstab(y_test, preds, rownames=['Actual'], colnames=['Predicted'])
fig, (ax1) = plt.subplots(ncols=1, figsize=(5,5))
sns.heatmap(cm,
```



```
In [74]: # The decision tree creation
tree.plot_tree(model7)
```

plt.savefig('DT.pdf')

plt.show()

Out[74]:

Out[74]:

```
In [75]: from sklearn.externals.six import StringIO
from IPython.display import Image
from sklearn.tree import export_graphviz
import pydotplus
from sklearn.tree import DecisionTreeRegressor
```

/Users/wangtingting/opt/anaconda3/lib/python3.7/site-packages/sklearn/externals/six.py:31: FutureWarning: the module is deprecated in version 0.21 and will be removed in version 0.23 since we've dropped support for Python 2.7. Please rely on the official version of six (<https://pypi.org/project/six/>).
*(<https://pypi.org/project/six/>),", FutureWarning)

```
In [76]: # build a tree graph
dot_data = StringIO()
export_graphviz(clf_entropy, out_file=dot_data,
                filled=True, rounded=True,
                special_characters=True)
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
Image(graph.create_png())
```



```
In [ ]:
```

```
In [ ]:
```