

```
from IPython.core.interactiveshell import InteractiveShell
InteractiveShell.ast_node_interactivity = 'all'
```

```
import warnings
warnings.filterwarnings('ignore')
warnings.simplefilter('ignore')
```

```
import sys
```

```
!load_ext pycodestyle_magic
```

```
pycodestyle.on
```

```
pycodestyle.off
```

```
import pandas as pd
import numpy as np
```

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
from pandas_profiling import ProfileReport
from pandas.plotting import scatter_matrix
```

```
!matplotlib inline
```

## index

- 1. Load the data
- 2. Review the data
  - 2.1 missing value and basic description
  - 2.2 data profile
  - 2.3 Analysis the features
- 3. Modeling
  - 3.1 training the model
  - 3.2 feature importance and explainability
- 4. Conclusion

## 1. Load the data

```
input_path = "default_of_credit_card_clients.xls"
df = pd.read_excel(input_path, header=0)
df.info()
df.head()
df.shape
```

```
Out[57]:
```

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	2	34	0	0	0	0	...	14331	14948	15549
3	4	50000	2	2	1	37	0	0	0	0	0	...	28314	28959	29547
4	5	50000	1	2	1	57	-1	0	-1	0	...	20940	19146	19131	...

5 rows x 25 columns

```
df.info()
```

(30000, 25)

## 2.Review the data

### missing value & the basic description

```
df.isna().sum()
```

ID	0
LIMIT_BAL	0
SEX	0
EDUCATION	0
MARRIAGE	0
AGE	0
PAY_0	0
PAY_2	0
PAY_3	0

```
5 rows x 25 columns
```

```
Out[57]:
```

```
(30000, 25)
```

## 2. Review the data

### missing value & the basic description

```
In [22]:
```

```
df.isna().sum()
```

```
Out[22]:
```

```
ID 0
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
dtype: int64
```

There is no missing value in all the columns , so we don't need to deal with the missing value.

```
In [17]:
```

```
df.describe().T
```

	ID	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	...	BILL_AMT1	BILL_AMT5	BILL_AMT6	PAY_A
count	30000.0	30000.0	30000.0	30000.0	30000.0	30000.0	30000.0	30000.0	30000.0	30000.0	...	30000.0	30000.0	30000.0	30000.0
mean	16784.322667	16784.322667	1.603733	1.853133	1.551867	35.485500	-0.167600	-0.166200	-0.220667	-0.266200	...	5921.163500	5921.163500	5921.163500	5921.163500
std	8660.398374	8660.398374	0.489129	0.790349	0.521970	9.217904	1.123802	1.168668	1.169139	1.133187	...	7363.860576	7363.860576	7363.860576	7363.860576
min	1	1	1	1	1	1	-2	-2	-2	-2	...	-165580	-165580	-165580	-165580
25%	15000.0	15000.0	1.000000	1.000000	1.000000	28.000000	-1.000000	-1.000000	-1.000000	-1.000000	...	5000.000000	5000.000000	5000.000000	5000.000000
50%	16784.322667	16784.322667	1.000000	1.000000	1.000000	35.485500	-0.167600	-0.166200	-0.220667	-0.266200	...	5921.163500	5921.163500	5921.163500	5921.163500
75%	22500.25	22500.25	2.000000	2.000000	2.000000	41.000000	0.000000	0.000000	0.000000	0.000000	...	67991.00	67991.00	67991.00	67991.00
max	30000.00	30000.00	2.000000	2.000000	2.000000	79.000000	8.000000	8.000000	8.000000	8.000000	...	964511.0	964511.0	964511.0	964511.0

```
In [18]:
```

```
df.dtypes
```

```
ID int64
LIMIT_BAL int64
SEX int64
EDUCATION int64
MARRIAGE int64
AGE int64
PAY_0 int64
PAY_2 int64
PAY_3 int64
PAY_4 int64
PAY_5 int64
PAY_6 int64
BILL_AMT1 int64
BILL_AMT2 int64
BILL_AMT3 int64
BILL_AMT4 int64
BILL_AMT5 int64
BILL_AMT6 int64
PAY_AMT1 int64
PAY_AMT2 int64
PAY_AMT3 int64
PAY_AMT4 int64
PAY_AMT5 int64
PAY_AMT6 int64
default payment next month int64
dtype: object
```

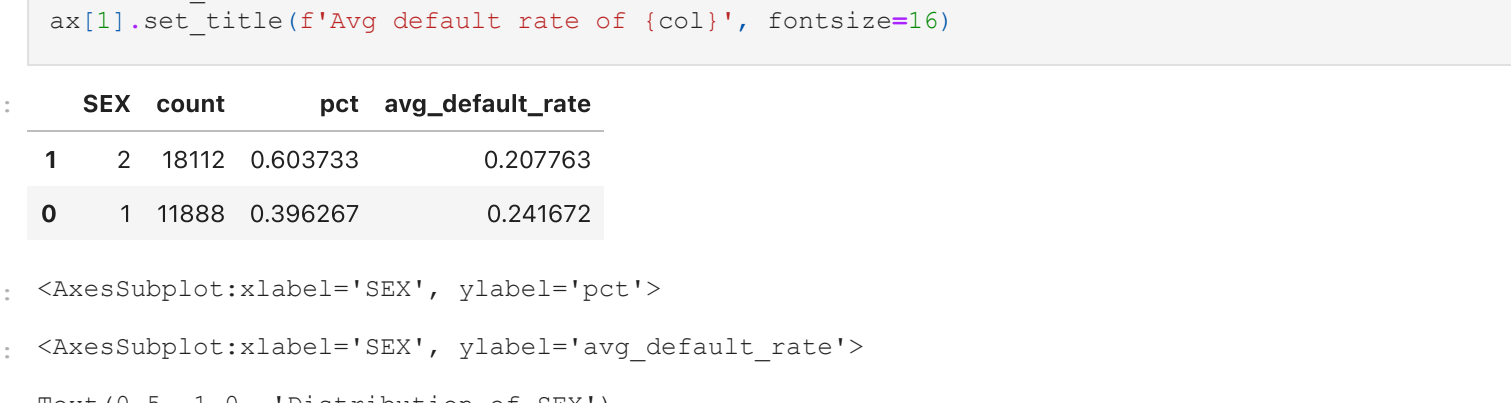
## check the outliers

```
In [183]:
```

```
sns.set(rc={'figure.figsize':(11.7,8.27)})
sns.boxplot(data=df.drop(columns='ID'))
```

```
Out[183]:
```

```
<AxesSubplot>
```



BILL\_AMTOUT and PAY\_AMTOUT seems have some outliers in a large positive value. But because it the Money amount, it is likely to be a long tail distribution. So I think it make sense, I don't deal with these "outliers".

## data-profile

```
In [21]:
```

```
from pandas_profiling import ProfileReport
profile = ProfileReport(df, title="Report_ModelingTest")
profile.to_file("Report_ModelingTest.html")
```

from the html files, we have some feeling about the data.

- The age variable has several high histograms.
- The BILL\_AMT1-6 contains some negative values.
- The default rate is 22% percent.
- The pay\_1 ~ pay\_6 are highly correlated with bill\_amt1 ~ bill\_amt6
- The bill\_amt1 ~ bill\_amt6 are highly correlated with pay\_amt1 ~ pay\_amt6

Next, we will check the above information and get a deeper understanding of the data.

## Analysis the features

### 1. Age

```
In [63]:
```

```
fig,ax = plt.subplots(nrows=1,ncols=2,figsize=(20,8))
col = 'AGE'
sns.histplot(data=df,x=col,ax=ax[0])
ax[0].set_title('Distribution of col', fontsize=16)
sns.barplot(data=df,x=col,y='default payment next month',ax=ax[1])
ax[1].set_title('Avg default rate of col', fontsize=16)
```

```
Out[63]:
```

```
<AxesSubplot: xlabel='AGE', ylabel='Count'>
```

```
Out[63]:
```

```
Text(0.5, 1.0, 'Distribution of AGE')
```

```
Out[63]:
```

```
<AxesSubplot: xlabel='AGE', ylabel='default payment next month'>
```

```
Out[63]:
```

```
Text(0.5, 1.0, 'Avg default rate of AGE')
```



We can see that the clients whose age are lower<24> and higher(>72) may have a higher default rate.

### 2. SEX

```
In [64]:
```

```
fig,ax = plt.subplots(nrows=1,ncols=2,figsize=(20,8))
data = df
col = 'SEX'
target = 'default payment next month'
count = data.groupby(col)[target].count()
pct = count/data.shape[0]
rate = data.groupby(col)[target].mean()
tmp = pd.concat([count,pct,rate],axis=1)
tmp.columns = ['count','pct','avg_default_rate']
tmp=tmp.reset_index(1).sort_values(by='avg_default_rate')
tmp[col] = tmp[col].apply(lambda x: x if x>0 else 0)
```

```
Out[64]:
```

	SEX	count	pct	avg_default_rate
0	1	18112	0.603733	0.207763
1	2	11888	0.396267	0.241672

```
Out[64]:
```

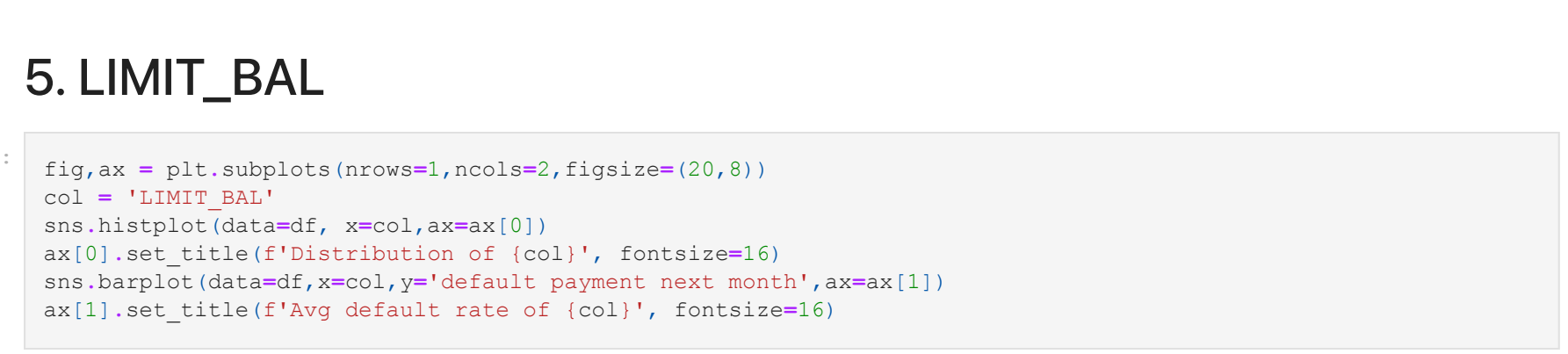
```
<AxesSubplot: xlabel='SEX', ylabel='pct'>
```

```
Out[64]:
```

```
Text(0.5, 1.0, 'Distribution of SEX')
```

```
Out[64]:
```

```
Text(0.5, 1.0, 'Avg default rate of SEX')
```



Male i.e. SEX==1 has 40% in the sample, but 25% default rate which is higher than female.

### 3. EDUCATION

```
In [73]:
```

```
fig,ax = plt.subplots(nrows=1,ncols=2,figsize=(20,8))
data = df
col = 'EDUCATION'
target = 'default payment next month'
d = [0:'zero',1:'graduate school',2:'university',3:'high school',4:'others',5:'unknown1',6:'unknown2']
count = data.groupby(col)[target].count()
pct = count/data.shape[0]
rate = data.groupby(col)[target].mean()
tmp = pd.concat([count,pct,rate],axis=1)
tmp.columns = ['count','pct','avg_default_rate']
tmp=tmp.reset_index(1).sort_values(by='avg_default_rate')
tmp[col] = tmp[col].apply(lambda x: x if x>0 else 0)
```

```
Out[73]:
```

	EDUCATION	count	pct	avg_default_rate
0	zero	14	0.000467	0.000000
1	graduate school	10585	0.352833	0.192348
2	university	14030	0.467667	0.237349
3	high school	4917	0.163900	0.291576
4	others	123	0.004100	0.056911
5	unknown1	280	0.009333	0.064286
6	unknown2	51	0.001700	0.156863

```
Out[73]:
```

```
<AxesSubplot: xlabel='EDUCATION', ylabel='pct'>
```

```
Out[73]:
```

```
<AxesSubplot: xlabel='EDUCATION', ylabel='avg_default_rate'>
```

```
Out[73]:
```

```
Text(0.5, 1.0, 'Distribution of EDUCATION')
```

```
Out[73]:
```

```
Text(0.5, 1.0, 'Avg default rate of EDUCATION')
```



We have some unknown values in the Education columns, like value 0,5,6. But this unknown value only have about 1% in the data and the default rate is quite low. We can just ignore these value. It won't affect the result.

But from the knowing value 1~3, we can find out that the higher the education level, the lower the risk.

If we want to get more meaningful variable, we should reorder this columns by using the average of the target like below, but in this case, I don't use the reorder.

```
Out [60]:
```

```
data = df
col = 'EDUCATION'
target = 'default payment next month'
rate = data.groupby(col)[target].mean()
rate = rate.sort_values(0)
rate
```

```
Out[60]:
```

```
reorder_edu = [k,v for k,v in zip(rate.index, range(len(rate)))]
reorder_edu
```

```
Out[60]:
```

```
df['EDUCATION2'] = df['EDUCATION'].apply(lambda x: reorder_edu[x])
df[['EDUCATION2']].sample(5)
```

```
Out[60]:
```

	EDUCATION	EDUCATION2
0	zero	3
2556	2	5
231	2	5
22459	2	5
24650	2	5

### 4. MARRIAGE

```
In [83]:
```

```
fig,ax = plt.subplots(nrows=1,ncols=2,figsize=(20,8))
data = df
col = 'MARRIAGE'
target = 'default payment next month'
d = [0:'zero',1:'married',2:'single',3:'others']
count = data.groupby(col)[target].count()
pct = count/data.shape[0]
rate = data.groupby(col)[target].mean()
tmp = pd.concat([count,pct,rate],axis=1)
tmp.columns = ['count','pct','avg_default_rate']
tmp=tmp.reset_index(1).sort_values(by='avg_default_rate')
tmp[col] = tmp[col].apply(lambda x: x if x>0 else 0)
```

```
Out[83]:
```

	MARRIAGE	count	pct	avg_default_rate
0	zero	54	0.001800	0.092593
1	married	13659	0.455300	0.234717
2	single	15964	0.532133	0.209283
3	others	323	0.010767	0.260062

```
Out[83]:
```

```
<AxesSubplot: xlabel='MARRIAGE', ylabel='pct'>
```

```
Out[83]:
```

```
<AxesSubplot: xlabel='MARRIAGE', ylabel='avg_default_rate'>
```

```
Out[83]:
```

```
Text(0.5, 1.0, 'Distribution of MARRIAGE')
```

```
Out[83]:
```

```
<AxesSubplot: xlabel='MARRIAGE', ylabel='avg_default_rate'>
```

```
Out[83]:
```

```
Text(0.5, 1.0, 'Distribution of MARRIAGE')
```

```
Out[83]:
```

```
Text(0.5, 1.0, 'Avg default rate of MARRIAGE')
```



"Others" marriage status has a higher risk than the other two. Married status has a higher risk than single.

We has some unknown value in the MARRIAGE, like the 0. But because the percentage of it is quite low, so I don't deal with the zeros.

### 5. LIMIT\_BAL

```
In [87]:
```

```
fig,ax = plt.subplots(nrows=1,ncols=2,figsize=(20,8))
col = 'LIMIT_BAL'
sns.histplot(data=df,x=col,ax=ax[0])
ax[0].set_title('Distribution of col', fontsize=16)
sns.barplot(data=df,x=col,y='default payment next month',ax=ax[1])
ax[1].set_title('Avg default rate of col', fontsize=16)
```

```
Out[87]:
```

```
<AxesSubplot: xlabel='LIMIT_BAL', ylabel='Count'>
```

```
Out[87]:
```

```
Text(0.5, 1.0, 'Distribution of LIMIT_BAL')
```

```
Out[87]:
```

```
<AxesSubplot: xlabel='LIMIT_BAL', ylabel='default payment next month'>
```

```
Out[87]:
```

```
Text(0.5, 1.0, 'Avg default rate of LIMIT_BAL')
```



The relationship between the LIMIT\_BAL and the average default rate is a U-shape.

### 6. PAY\_0

```
In [89]:
```

```
fig,ax = plt.subplots(nrows=1,ncols=2,figsize=(20,8))
col = 'PAY_0'
sns.countplot(data=df,x=col,ax=ax[0])
ax[0].set_title('Distribution of col', fontsize=16)
sns.barplot(data=df,x=col,y='default payment next month',ax=ax[1])
ax[1].set_title('Avg default rate of col', fontsize=16)
```

```
Out[89]:
```

```
<AxesSubplot: xlabel='PAY_0', ylabel='count'>
```

```
Out[89]:
```

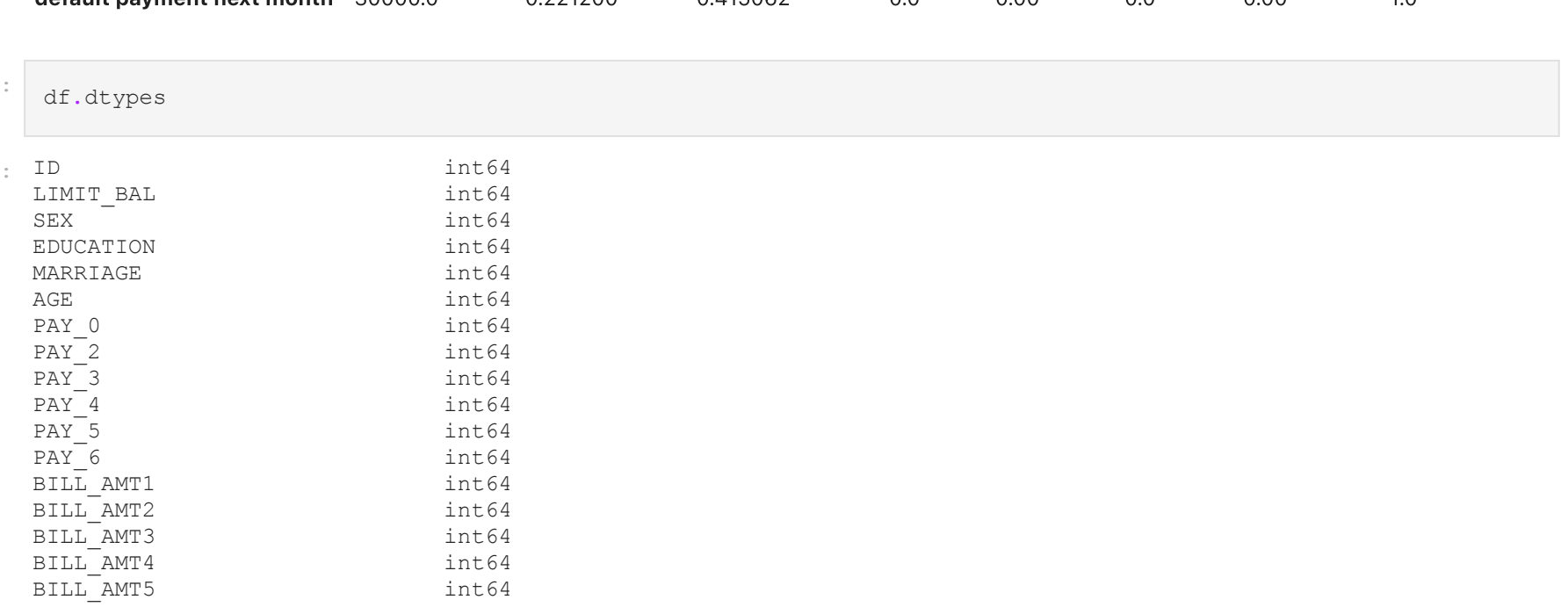
```
Text(0.5, 1.0, 'Distribution of PAY_0')
```

```
Out[89]:
```

```
<AxesSubplot: xlabel='PAY_0', ylabel='default payment next month'>
```

```
Out[89]:
```

```
Text(0.5, 1.0, 'Avg default rate of PAY_0')
```



In general, the higher the PAY\_0, the Higher risk.

### 7. BILL\_AMT1

```
In [140]:
```

```
fig,ax = plt.subplots(nrows=1,ncols=2,figsize=(20,8))
col = 'BILL_AMT1'
sns.histplot(data=df,x=col,ax=ax[0])
ax[0].set_title('Distribution of col', fontsize=16)
sns.distplot(df[df[target]==0][col],label='No Default',ax=ax[1])
sns.distplot(df[df[target]==1][col],label='Default',ax=ax[1])
plt.legend()
ax[1].set_title('Distribution of col in different target', fontsize=16)
```

```
Out[140]:
```

```
<AxesSubplot: xlabel='BILL_AMT1', ylabel='Count'>
```

```
Out[140]:
```

```
Text(0.5, 1.0, 'Distribution of BILL_AMT1')
```

```
Out[140]:
```

```
<AxesSubplot: xlabel='BILL_AMT1', ylabel='Density'>
```

```
Out[140]:
```

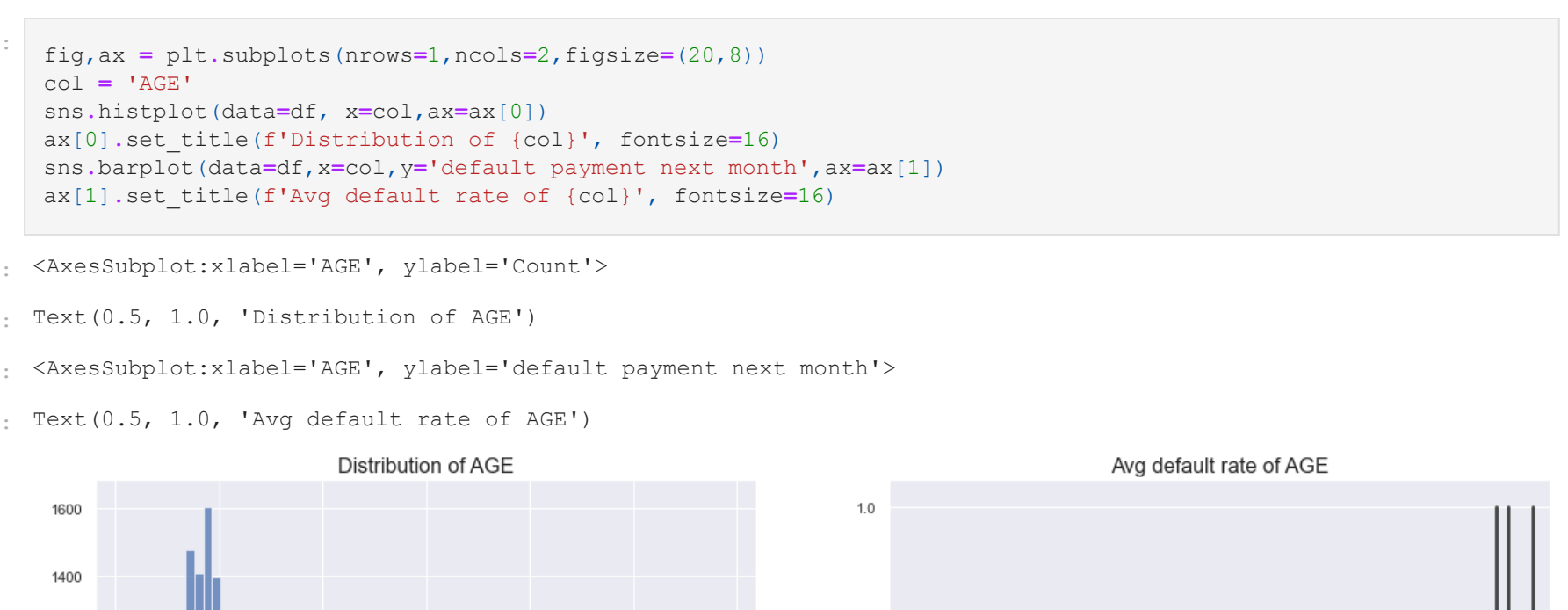
```
<AxesSubplot: xlabel='BILL_AMT1', ylabel='Density'>
```

```
Out[140]:
```

```
<matplotlib.legend.Legend at 0x7fe60a07190>
```

```
Out[140]:
```

```
Text(0.5, 1.0, 'Distribution of BILL_AMT1 in different target')
```



```
In [140]:
```

```
num_bins=20
q = [round(i*(1/num_bins),3) for i in range(num_bins)]
bins
```



