

```
In [43]: #outlier treatment
ldi["PRI.SANCTIONED.AMOUNT"].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])

Out[43]:
0.10      0.00
0.25      0.00
0.50      0.00
0.70      43486.20
0.90      447157.30
0.95      1030280.05
0.99      3467119.37
Name: PRI.SANCTIONED.AMOUNT, dtype: float64
```

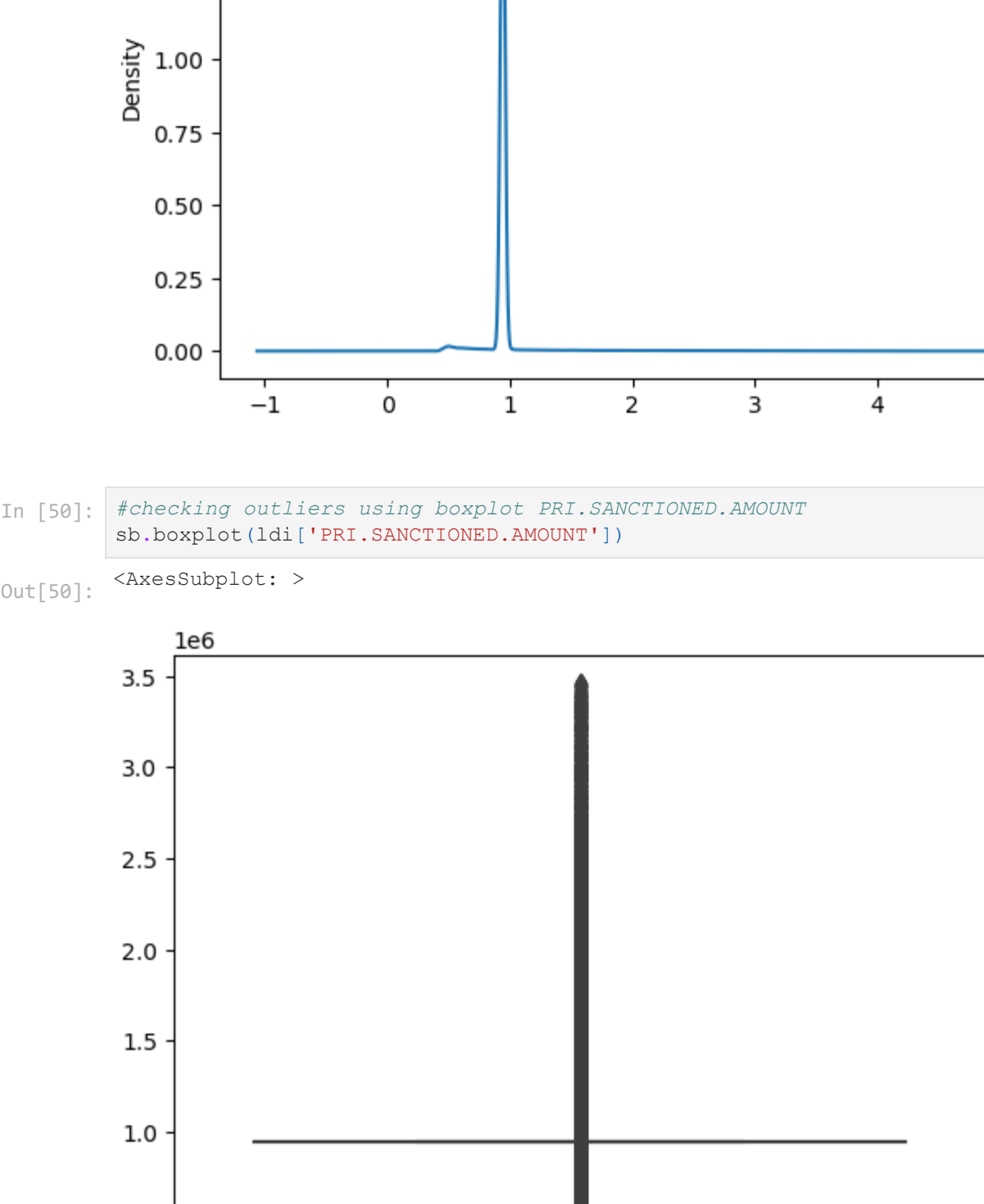
```
In [44]: #seeing limits
psa_hb=ldi[ldi["PRI.SANCTIONED.AMOUNT"] > 3467119.37].copy()
psa_lb=ldi[ldi["PRI.SANCTIONED.AMOUNT"] < 447157.30].copy()
```

```
In [45]: #putting them at one place
print(len(psa_hb))
print(len(psa_lb))
```

```
2332
209838
```

```
In [46]: #removing outliers
for x in iidi["PRI.SANCTIONED.AMOUNT"]:
    if x > 3467119.37:
        ldi["PRI.SANCTIONED.AMOUNT"].replace(x,np.nan,inplace=True)
for x in iidi["PRI.SANCTIONED.AMOUNT"]:
    if x < 447157.30:
        ldi["PRI.SANCTIONED.AMOUNT"].replace(x,np.nan,inplace=True)
```

```
In [47]: #checking outliers gone or not
sb.boxplot(ldi["PRI.SANCTIONED.AMOUNT"])
```



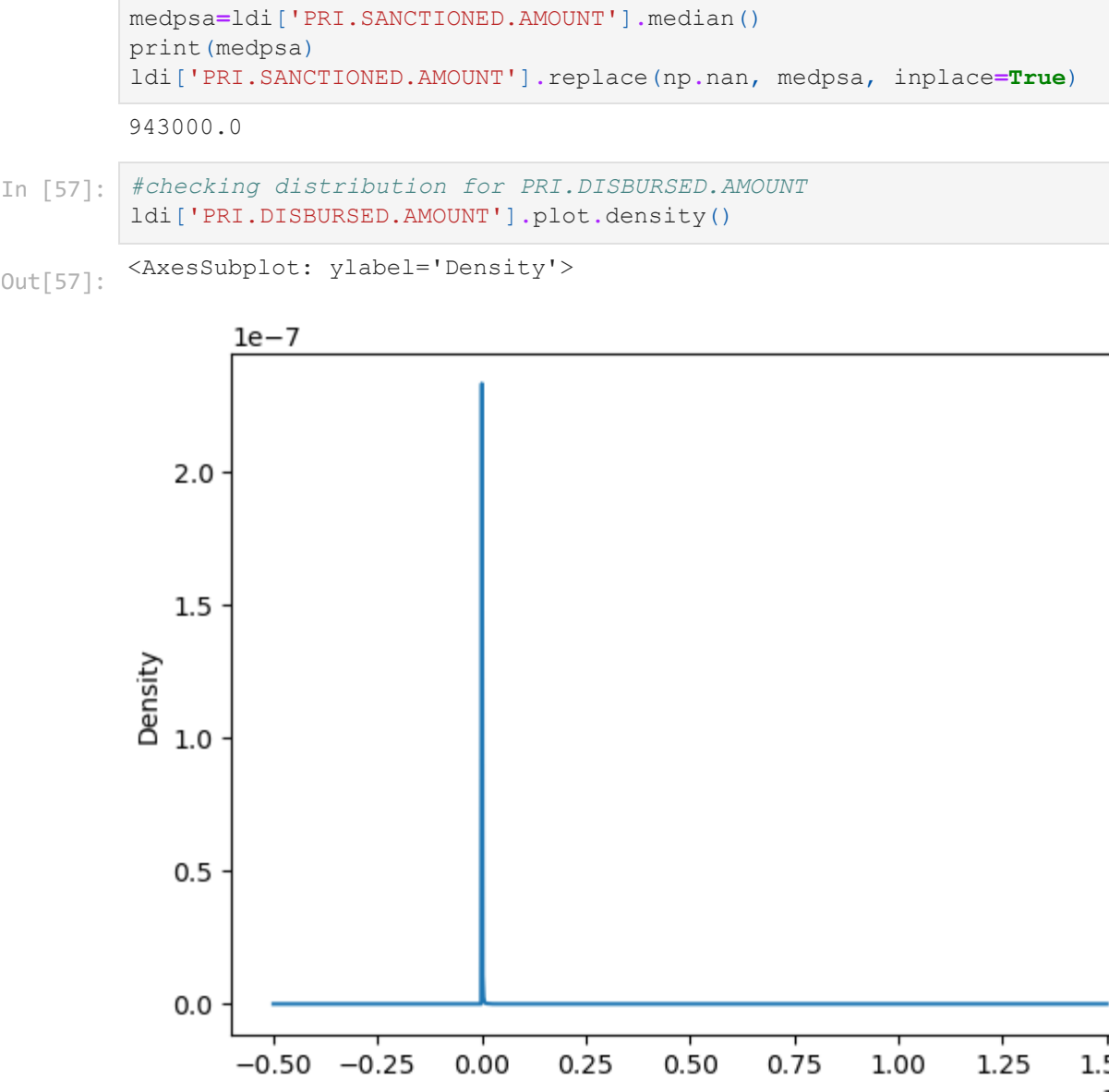
```
In [48]: #replacing with median values
medpsa=ldi["PRI.SANCTIONED.AMOUNT"].median()
print(medpsa)
ldi["PRI.SANCTIONED.AMOUNT"].replace(np.nan, medpsa, inplace=True)
```

```
In [49]: #checking distribution for PRI.SANCTIONED.AMOUNT
ldi["PRI.SANCTIONED.AMOUNT"].plot.density()
```

<AxesSubplot: ylabel='Density'>



```
In [50]: #checking outliers using boxplot PRI.DISBURSED.AMOUNT
sb.boxplot(ldi["PRI.DISBURSED.AMOUNT"])
```



```
In [51]: #outlier treatment
ldi["PRI.DISBURSED.AMOUNT"].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])

Out[51]:
0.10      943000.0
0.25      943000.0
0.50      943000.0
0.70      943000.0
0.90      943000.0
0.95      943000.0
0.99      2250000.0
Name: PRI.DISBURSED.AMOUNT, dtype: float64
```

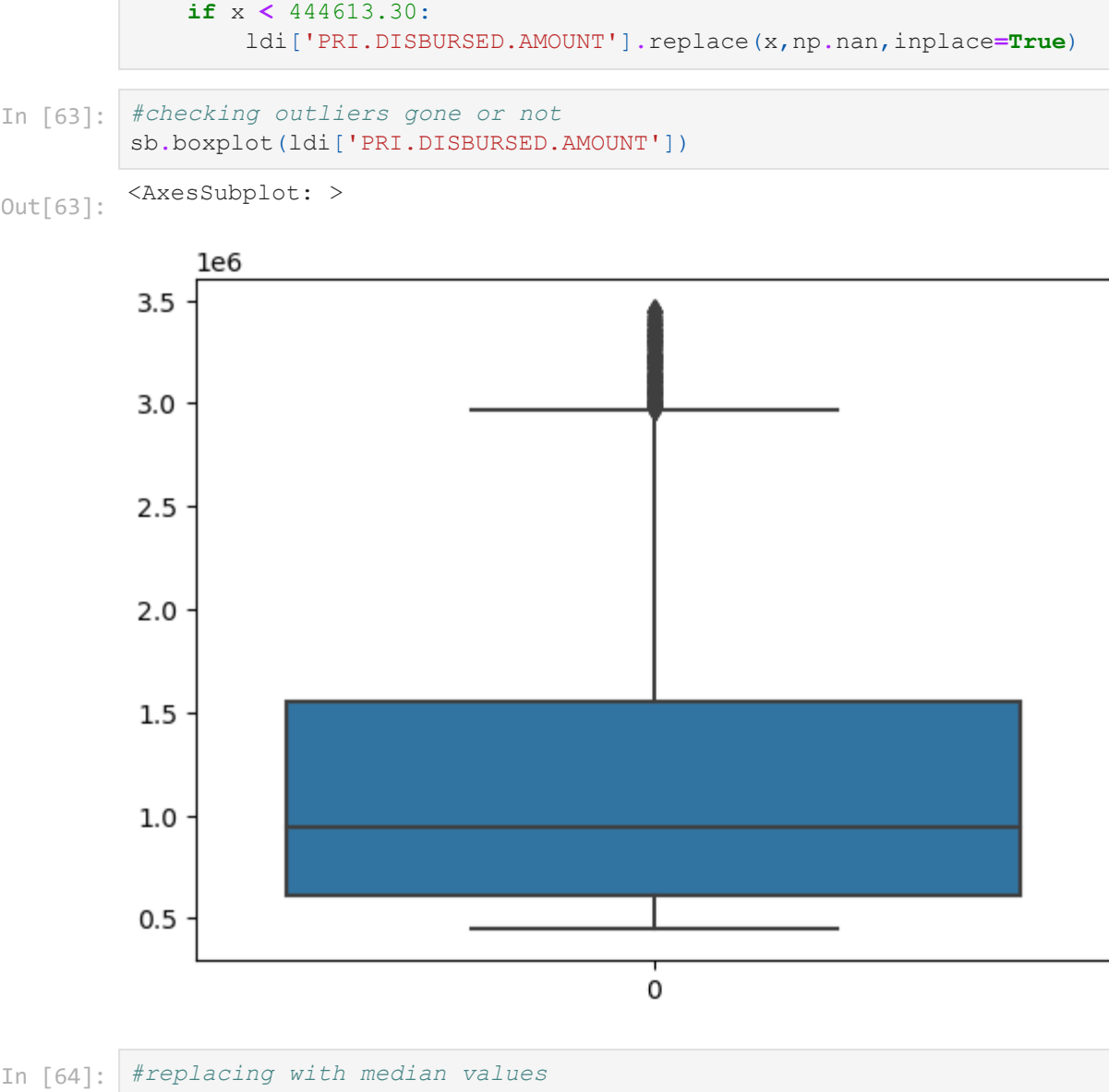
```
In [52]: #seeing limits
psa_hb=ldi[ldi["PRI.DISBURSED.AMOUNT"] > 2250000.0].copy()
psa_lb=ldi[ldi["PRI.DISBURSED.AMOUNT"] < 943000.0].copy()
```

```
In [53]: #putting them at one place
print(len(psa_hb))
print(len(psa_lb))
```

```
222
10488
```

```
In [54]: #removing outliers
for x in iidi["PRI.DISBURSED.AMOUNT"]:
    if x > 2250000.0:
        ldi["PRI.DISBURSED.AMOUNT"].replace(x,np.nan,inplace=True)
for x in iidi["PRI.DISBURSED.AMOUNT"]:
    if x < 943000.0:
        ldi["PRI.DISBURSED.AMOUNT"].replace(x,np.nan,inplace=True)
```

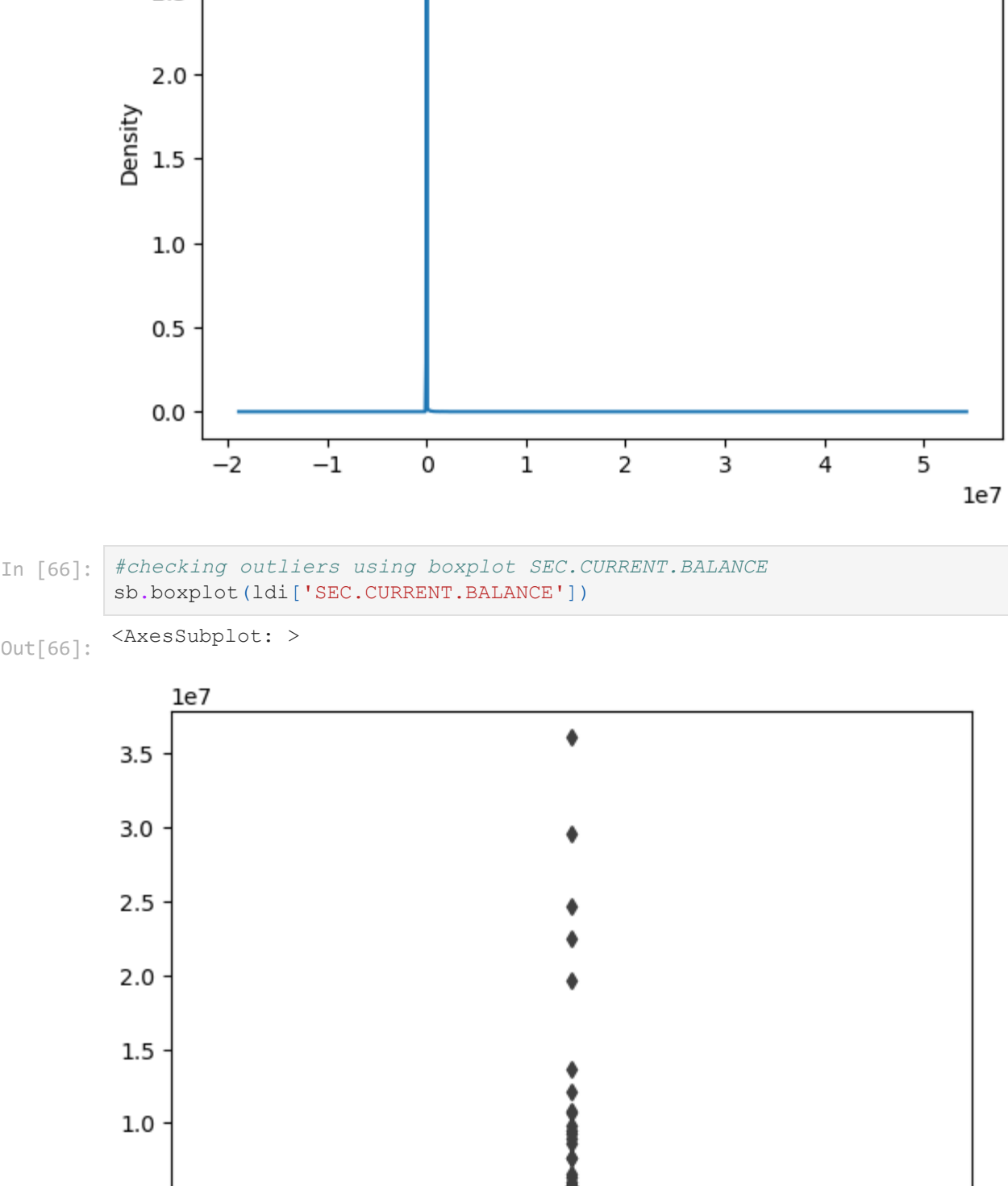
```
In [55]: #checking outliers gone or not
sb.boxplot(ldi["PRI.DISBURSED.AMOUNT"])
```



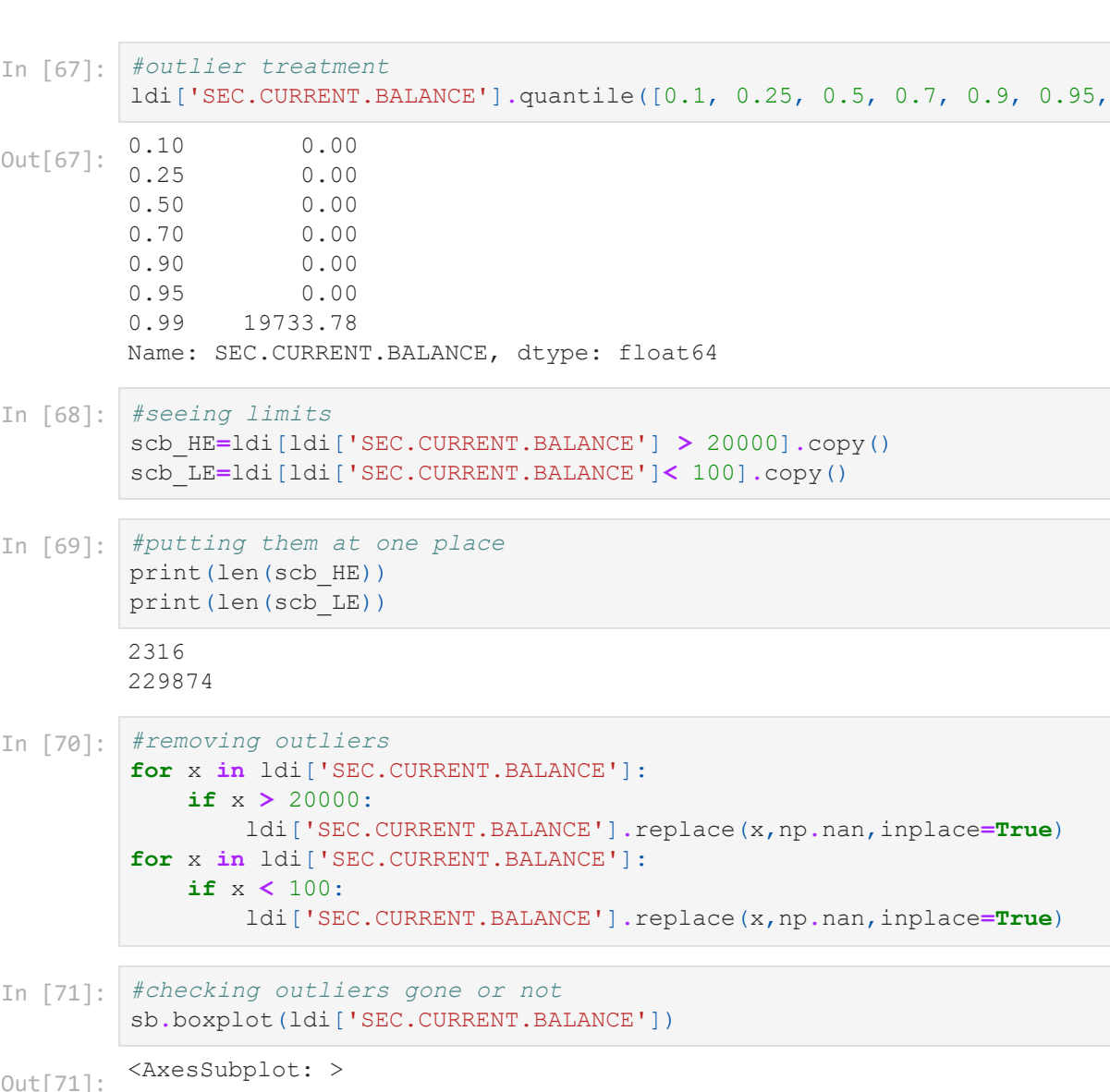
```
In [56]: #replacing with median values
medpsa=ldi["PRI.DISBURSED.AMOUNT"].median()
print(medpsa)
ldi["PRI.DISBURSED.AMOUNT"].replace(np.nan, medpsa, inplace=True)
```

```
In [57]: #checking distribution for PRI.DISBURSED.AMOUNT
ldi["PRI.DISBURSED.AMOUNT"].plot.density()
```

<AxesSubplot: ylabel='Density'>



```
In [58]: #checking outliers using boxplot PRI.DISBURSED.AMOUNT
sb.boxplot(ldi["PRI.DISBURSED.AMOUNT"])
```



```
In [59]: #outlier treatment
ldi["SEC.DISBURSED.AMOUNT"].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])

Out[59]:
0.10      0.00
0.25      0.00
0.50      0.00
0.70      40000.00
0.90      444613.30
0.95      1027992.85
0.99      3457573.88
Name: SEC.DISBURSED.AMOUNT, dtype: float64
```

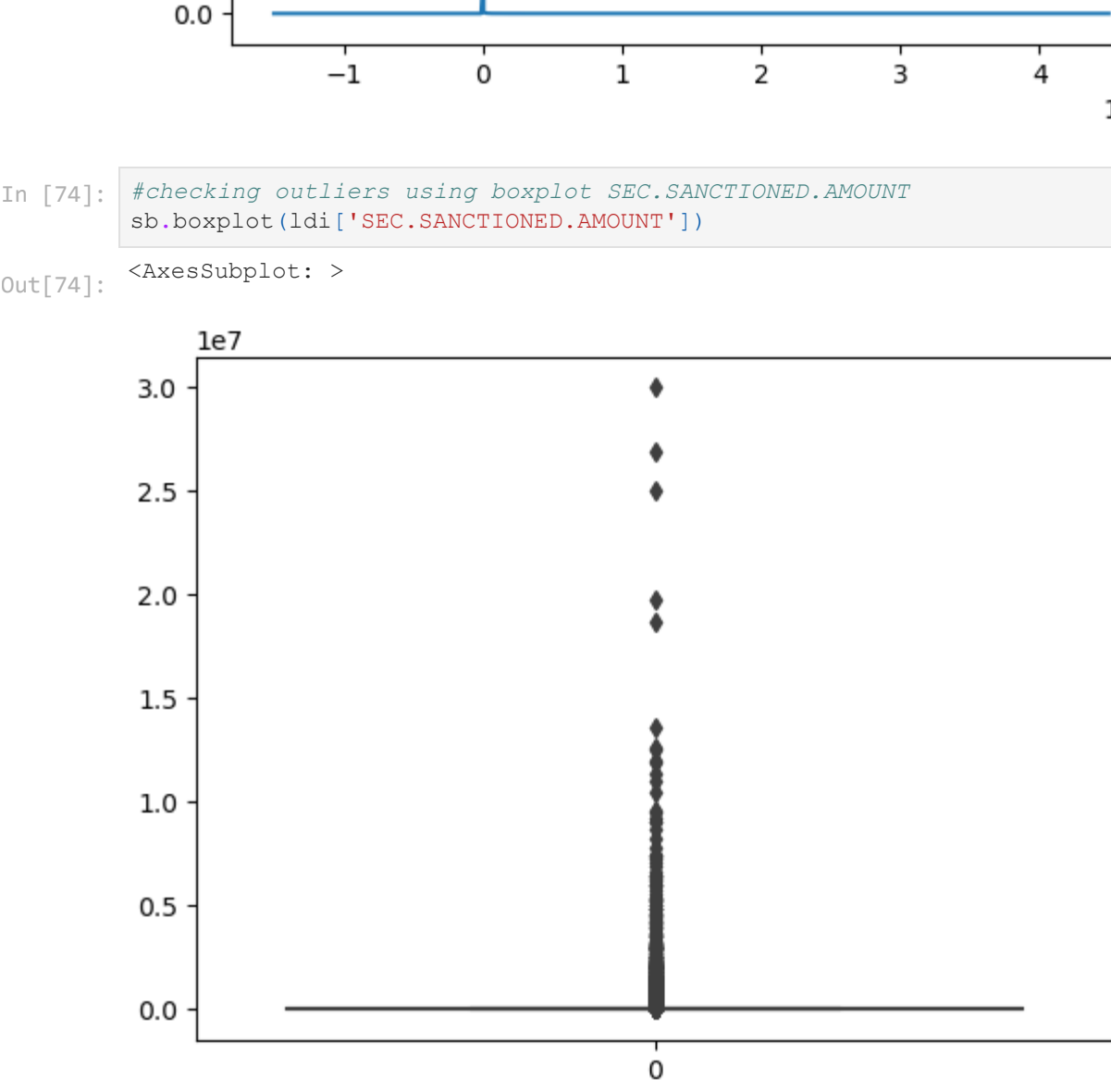
```
In [60]: #seeing limits
pda_hb=ldi[ldi["SEC.DISBURSED.AMOUNT"] > 3457573.88].copy()
pda_lb=ldi[ldi["SEC.DISBURSED.AMOUNT"] < 444613.30].copy()
```

```
In [61]: #putting them at one place
print(len(pda_hb))
print(len(pda_lb))
```

```
2332
209838
```

```
In [62]: #removing outliers
for x in iidi["SEC.DISBURSED.AMOUNT"]:
    if x > 3457573.88:
        ldi["SEC.DISBURSED.AMOUNT"].replace(x,np.nan,inplace=True)
for x in iidi["SEC.DISBURSED.AMOUNT"]:
    if x < 444613.30:
        ldi["SEC.DISBURSED.AMOUNT"].replace(x,np.nan,inplace=True)
```

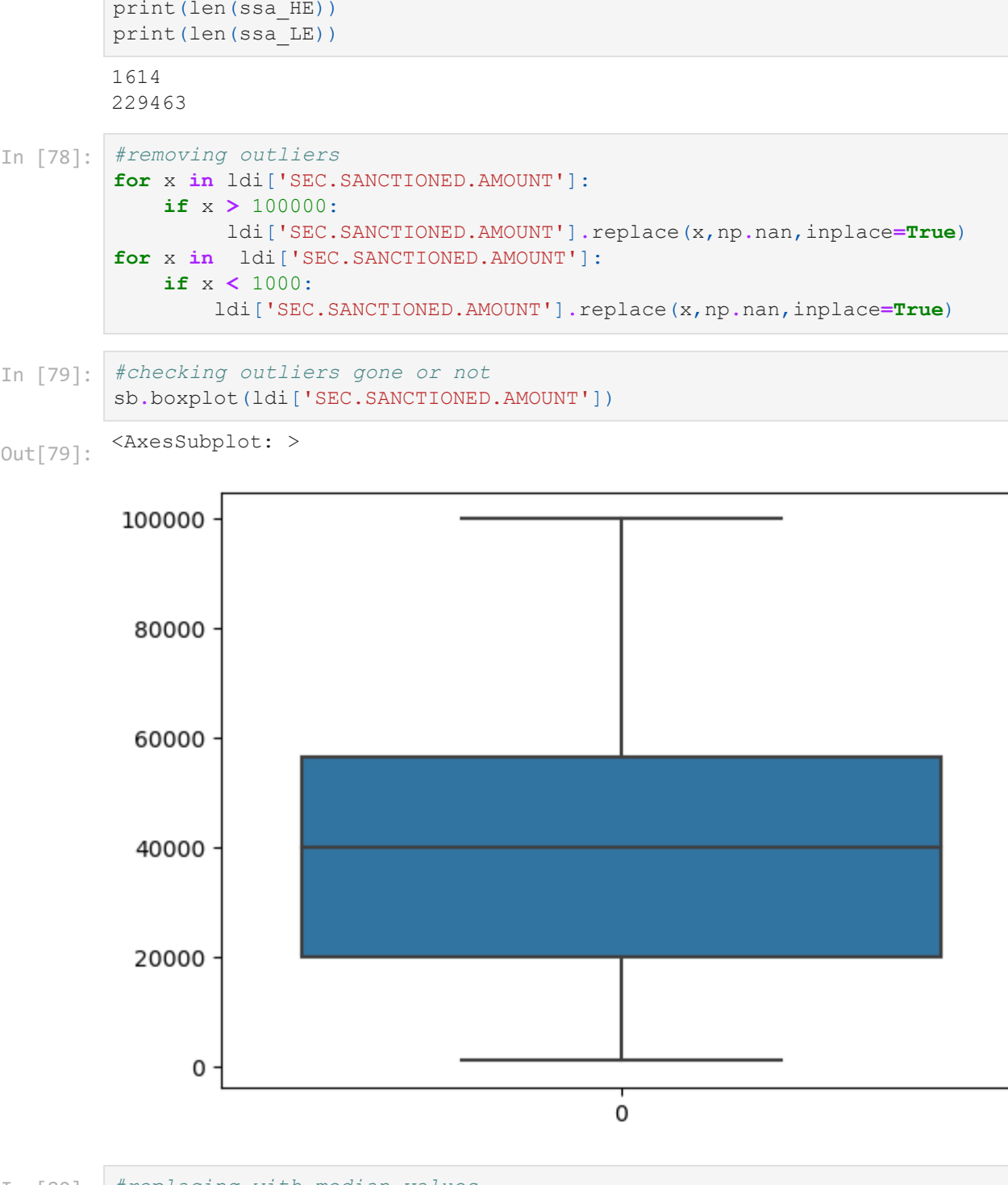
```
In [63]: #checking outliers gone or not
sb.boxplot(ldi["SEC.DISBURSED.AMOUNT"])
```



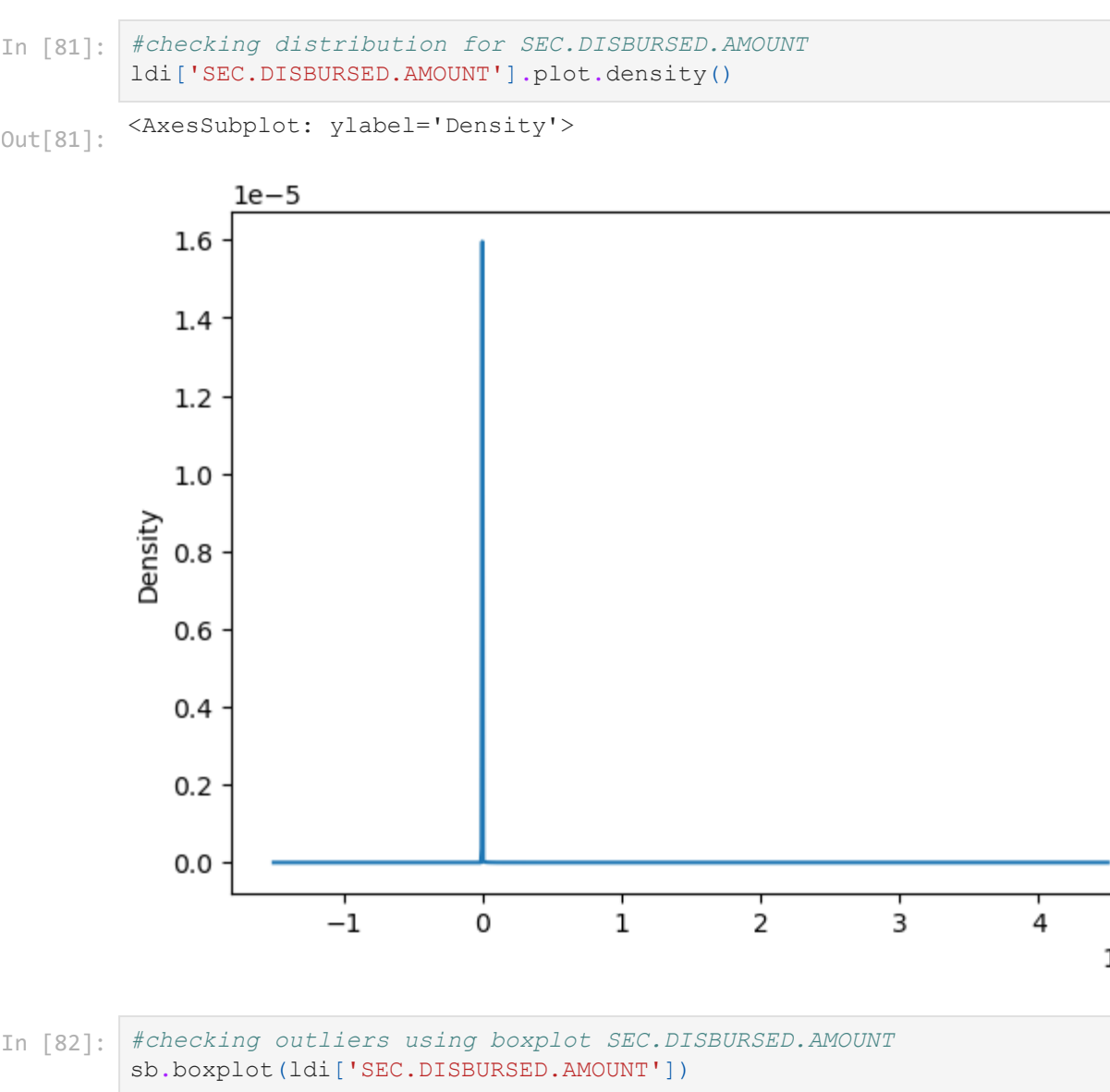
```
In [64]: #replacing with median values
medpda=ldi["SEC.DISBURSED.AMOUNT"].median()
print(medpda)
ldi["SEC.DISBURSED.AMOUNT"].replace(np.nan, medpda, inplace=True)
```

```
In [65]: #checking distribution for SEC.DISBURSED.AMOUNT
ldi["SEC.DISBURSED.AMOUNT"].plot.density()
```

<AxesSubplot: ylabel='Density'>



```
In [66]: #checking outliers using boxplot SEC.CURRENT.BALANCE
sb.boxplot(ldi["SEC.CURRENT.BALANCE"])
```



```
In [67]: #outlier treatment
ldi["SEC.CURRENT.BALANCE"].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])

Out[67]:
0.10      0.00
0.25      0.00
0.50      0.00
0.70      0.00
0.90      0.00
0.95      0.00
0.99      19733.78
Name: SEC.CURRENT.BALANCE, dtype: float64
```

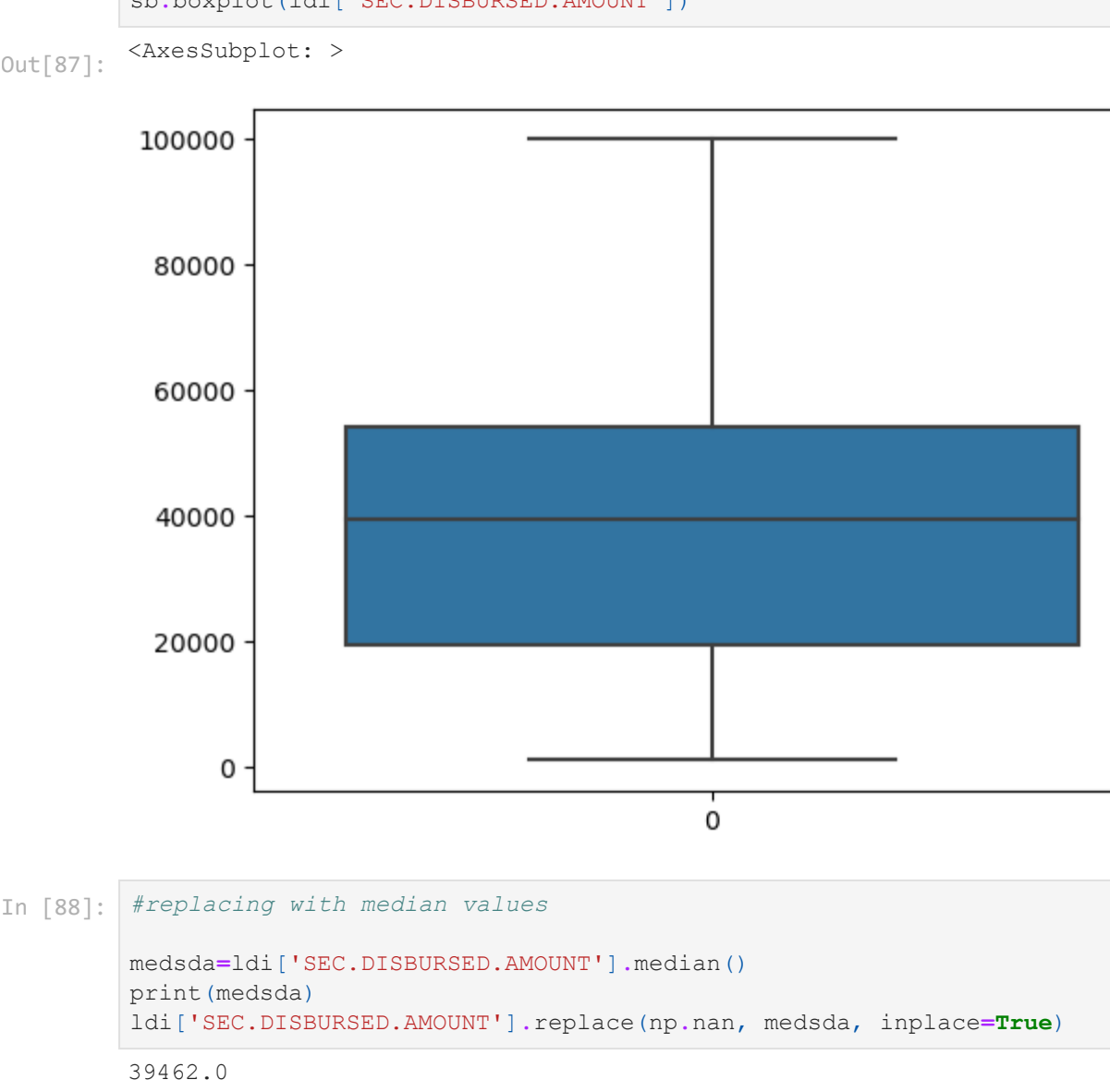
```
In [68]: #seeing limits
scb_hb=ldi[ldi["SEC.CURRENT.BALANCE"] > 20000].copy()
scb_lb=ldi[ldi["SEC.CURRENT.BALANCE"] < 100].copy()
```

```
In [69]: #putting them at one place
print(len(scb_hb))
print(len(scb_lb))
```

```
2316
229874
```

```
In [70]: #removing outliers
for x in iidi["SEC.CURRENT.BALANCE"]:
    if x > 20000:
        ldi["SEC.CURRENT.BALANCE"].replace(x,np.nan,inplace=True)
for x in iidi["SEC.CURRENT.BALANCE"]:
    if x < 100:
        ldi["SEC.CURRENT.BALANCE"].replace(x,np.nan,inplace=True)
```

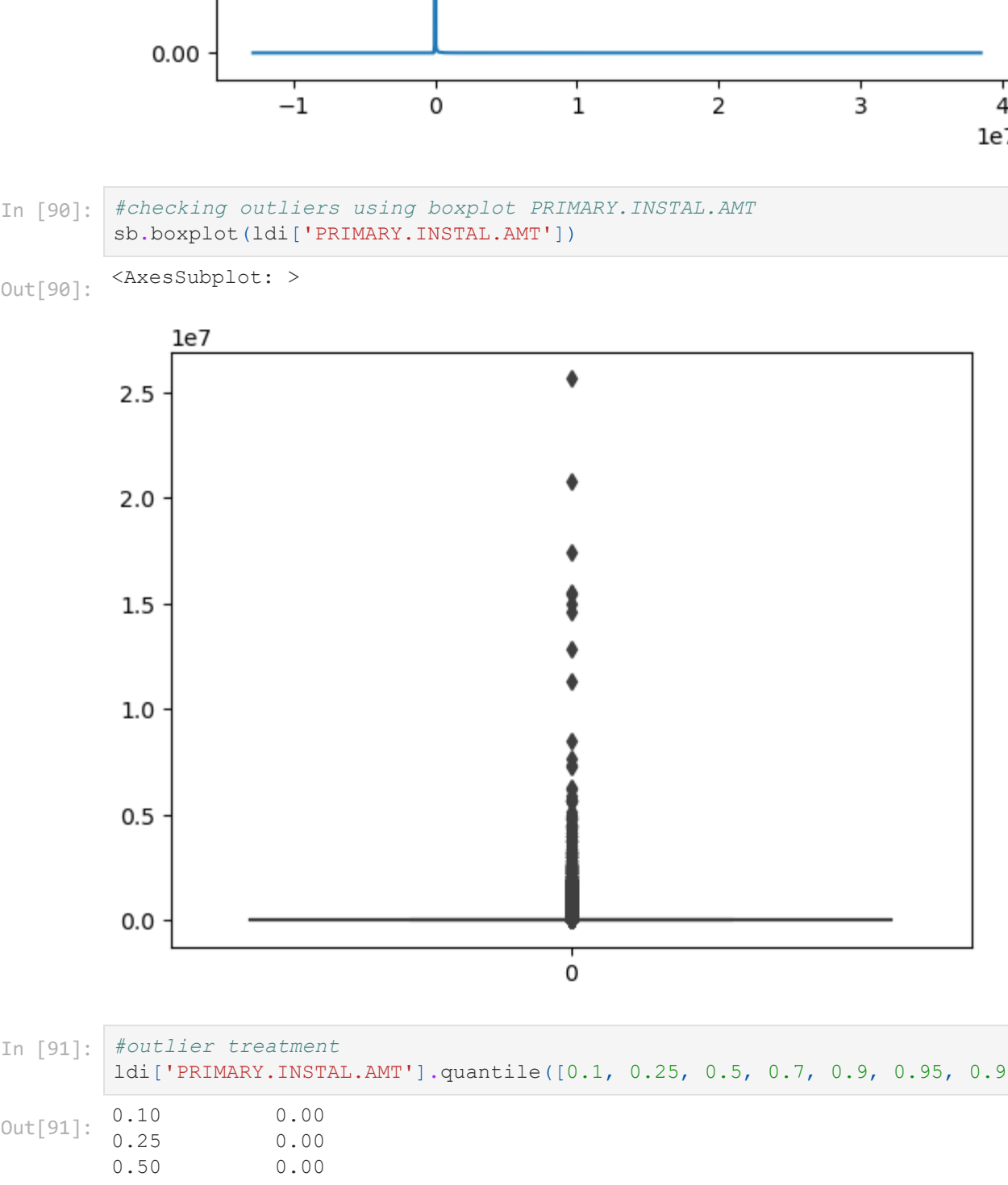
```
In [71]: #checking outliers gone or not
sb.boxplot(ldi["SEC.CURRENT.BALANCE"])
```



```
In [72]: #replacing with median values
medscb=ldi["SEC.CURRENT.BALANCE"].median()
print(medscb)
ldi["SEC.CURRENT.BALANCE"].replace(np.nan, medscb, inplace=True)
```

```
In [73]: #checking distribution for SEC.SANCTIONED.AMOUNT
ldi["SEC.SANCTIONED.AMOUNT"].plot.density()
```

<AxesSubplot: ylabel='Density'>



```
In [74]: #checking outliers using boxplot SEC.SANCTIONED.AMOUNT
sb.boxplot(ldi["SEC.SANCTIONED.AMOUNT"])
```



```
In [75]: #outlier treatment
ldi["SEC.SANCTIONED.AMOUNT"].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])

Out[75]:
0.10      0.0
0.25      0.0
0.50      0.0
0.70      0.0
0.90      0.0
0.95      50000.0
0.99      50000.0
Name: SEC.SANCTIONED.AMOUNT, dtype: float64
```

```
In [76]: #seeing limits
ssa_hb=ldi[ldi["SEC.SANCTIONED.AMOUNT"] > 100000].copy()
ssa_lb=ldi[ldi["SEC.SANCTIONED.AMOUNT"] < 1000].copy()
```

```
In [77]: #putting them at one place
print(len(ssa_hb))
print(len(ssa_lb))
```

```
1614
229463
```

```
In [78]: #removing outliers
for x in iidi["SEC.SANCTIONED.AMOUNT"]:
    if x > 100000:
        ldi["SEC.SANCTIONED.AMOUNT"].replace(x,np.nan,inplace=True)
for x in iidi["SEC.SANCTIONED.AMOUNT"]:
    if x < 1000:
        ldi["SEC.SANCTIONED.AMOUNT"].replace(x,np.nan,inplace=True)
```

```
In [79]: #checking outliers gone or not
sb.boxplot(ldi["SEC.SANCTIONED.AMOUNT"])
```



```
In [80]: #replacing with median values
medssa=ldi["SEC.SANCTIONED.AMOUNT"].median()
print(medssa)
ldi["SEC.SANCTIONED.AMOUNT"].replace(np.nan, medssa, inplace=True)
```

```
In [81]: #checking distribution for SEC.DISBURSED.AMOUNT
ldi["SEC.DISBURSED.AMOUNT"].plot.density()
```

<AxesSubplot: ylabel='Density'>


```
In [82]: #checking outliers using boxplot SEC.DISBURSED.AMOUNT
sb.boxplot(ldi["SEC.DISBURSED.AMOUNT"])
```



```
In [83]: #outlier treatment
ldi["SEC.DISBURSED.AMOUNT"].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])

Out[83]:
0.10      0.00
0.25      0.00
0.50      0.00
0.70      0.00
0.90      0.00
0.95      0.00
0.99      47944.83
Name: SEC.DISBURSED.AMOUNT, dtype: float64
```

```
In [84]: #seeing limits
sda_hb=ldi[ldi["SEC.DISBURSED.AMOUNT"] > 100000].copy()
sda_lb=ldi[ldi["SEC.DISBURSED.AMOUNT"] < 1000].copy()
```

```
In [85]: #putting them at one place
print(len(sda_hb))
print(len(sda_lb))
```

```
1559
229516
```

```
In [86]: #removing outliers
for x in iidi["SEC.DISBURSED.AMOUNT"]:
    if x > 100000:
        ldi["SEC.DISBURSED.AMOUNT"].replace(x,np.nan,inplace=True)
for x in iidi["SEC.DISBURSED.AMOUNT"]:
    if x < 1000:
        ldi["SEC.DISBURSED.AMOUNT"].replace(x,np.nan,inplace=True)
```

```
In [87]: #checking outliers gone or not
sb.boxplot(ldi["SEC.DISBURSED.AMOUNT"])
```



```
In [88]: #replacing with median values
medsda=ldi["SEC.DISBURSED.AMOUNT"].median()
print(medsda)
ldi["SEC.DISBURSED.AMOUNT"].replace(np.nan, medsda, inplace=True)
```

```
In [89]: #checking distribution for PRIMARY.INSTAL.AMT
ldi["PRIMARY.INSTAL.AMT"].plot.density()
```

<AxesSubplot: ylabel='Density'>


```
In [90]: #checking outliers using boxplot PRIMARY.INSTAL.AMT
sb.boxplot(ldi["PRIMARY.INSTAL.AMT"])
```



```
In [91]: #outlier treatment
ldi["PRIMARY.INSTAL.AMT"].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])

Out[91]:
0.10      0.00
0.25      0.00
0.50      0.00
0.70      1156.00
0.90      11504.00
0.95      26376.45
0.99      256159.38
Name: PRIMARY.INSTAL.AMT, dtype: float64
```

```
In [92]: #seeing limits
pia_hb=ldi[ldi["PRIMARY.INSTAL.AMT"] > 15000].copy()
pia_lb=ldi[ldi["PRIMARY.INSTAL.AMT"] < 1200.00].copy()
```

```
In [93]: #putting them at one place
print(len(pia_hb))
print(len(pia_lb))
```

```
18672
163608
```

```
In [94]: #removing outliers
for x in iidi["PRIMARY.INSTAL.AMT"]:
    if x > 15000:
        ldi["PRIMARY.INSTAL.AMT"].replace(x,np.nan,inplace=True)
for x in iidi["PRIMARY.INSTAL.AMT"]:
    if x < 1200.00:
        ldi["PRIMARY.INSTAL.AMT"].replace(x,np.nan,inplace=True)
```

```
In [95]: #checking outliers gone or not
sb.boxplot(ldi["PRIMARY.INSTAL.AMT"])
```


In [95]:
<AxesSubplot: >

In [96]:
#replacing with median values
medialb=ldi['PRIMARY.INSTAL.AMT'].median()
print(medialb)
ldi['PRIMARY.INSTAL.AMT'].replace(np.nan, medialb, inplace=True)
3786.0
In [97]:
#checking distribution for SEC.INSTAL.AMT
ldi['SEC.INSTAL.AMT'].plot.density()
<AxesSubplot: ylabel='Density'>

In [98]:
#checking outliers using boxplot SEC.INSTAL.AMT
sb.boxplot(ldi['SEC.INSTAL.AMT'])
<AxesSubplot: >

In [99]:
#outlier treatment
ldi['SEC.INSTAL.AMT'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
Out[99]:
0.10 0.0
0.25 0.0
0.50 0.0
0.70 0.0
0.90 0.0
0.95 0.0
0.99 0.0
Name: SEC.INSTAL.AMT, dtype: float64
In [100]:
#seeing limits
sia_HB=ldi[ldi['SEC.INSTAL.AMT'] > 8000.00].copy()
sia_LB=ldi[ldi['SEC.INSTAL.AMT'] < 100.00].copy()
In [101]:
#putting them at one place
print(len(sia_HB))
print(len(sia_LB))
113
230974
In [102]:
#removing outliers
for x in ian_idi['SEC.INSTAL.AMT']:
 if x > 8000.00:
 ldi['SEC.INSTAL.AMT'].replace(x,np.nan,inplace=True)
 if x < 100.00:
 ldi['SEC.INSTAL.AMT'].replace(x,np.nan,inplace=True)
In [103]:
#checking outliers gone or not
sb.boxplot(ldi['SEC.INSTAL.AMT'])
<AxesSubplot: >

In [104]:
#replacing with median values
medialb=ldi['SEC.INSTAL.AMT'].median()
print(medialb)
ldi['SEC.INSTAL.AMT'].replace(np.nan, medialb, inplace=True)
2440.0
In [105]:
#checking distribution for Acc_age_months
ldi['Acc_age_months'].plot.density()
<AxesSubplot: ylabel='Density'>

In [106]:
#checking outliers using boxplot disbursed_amount
sb.boxplot(ldi['Acc_age_months'])
<AxesSubplot: >

In [107]:
#outlier treatment
ldi['Acc_age_months'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
Out[107]:
0.10 0.0
0.25 0.0
0.50 0.0
0.70 11.0
0.90 26.0
0.95 35.0
0.99 66.0
Name: Acc_age_months, dtype: float64
In [108]:
#seeing limits
aam_HB=ldi[ldi['Acc_age_months'] > 35].copy()
aam_LB=ldi[ldi['Acc_age_months'] < 11].copy()
In [109]:
#putting them at one place
print(len(aam_HB))
print(len(aam_LB))
1309
163133
In [110]:
#removing outliers
for x in ian_idi['Acc_age_months']:
 if x > 35:
 ldi['Acc_age_months'].replace(x,np.nan,inplace=True)
 if x < 11:
 ldi['Acc_age_months'].replace(x,np.nan,inplace=True)
In [111]:
#checking outliers gone or not
sb.boxplot(ldi['Acc_age_months'])
<AxesSubplot: >

In [112]:
#replacing with median values
medaam=ldi['Acc_age_months'].median()
print(medam)
ldi['Acc_age_months'].replace(np.nan, medaam, inplace=True)
17.0
In [113]:
#checking distribution for NO.OF_INQUIRIES
ldi['NO.OF_INQUIRIES'].plot.density()
<AxesSubplot: ylabel='Density'>

In [114]:
#checking outliers using boxplot NO.OF_INQUIRIES
sb.boxplot(ldi['NO.OF_INQUIRIES'])
<AxesSubplot: >

In [115]:
#outlier treatment
ldi['NO.OF_INQUIRIES'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
Out[115]:
0.10 0.0
0.25 0.0
0.50 0.0
0.70 0.0
0.90 1.0
0.95 1.0
0.99 1.0
Name: NO.OF_INQUIRIES, dtype: float64
In [116]:
#seeing limits
noi_HB=ldi[ldi['NO.OF_INQUIRIES'] > 5].copy()
noi_LB=ldi[ldi['NO.OF_INQUIRIES'] < 1].copy()
In [117]:
#putting them at one place
print(len(noi_HB))
print(len(noi_LB))
629
201961
In [118]:
#removing outliers
for x in ian_idi['NO.OF_INQUIRIES']:
 if x > 5:
 ldi['NO.OF_INQUIRIES'].replace(x,np.nan,inplace=True)
 if x < 1:
 ldi['NO.OF_INQUIRIES'].replace(x,np.nan,inplace=True)
In [119]:
#checking outliers gone or not
sb.boxplot(ldi['NO.OF_INQUIRIES'])
<AxesSubplot: >

In [120]:
#replacing with median values
mednoi=ldi['NO.OF_INQUIRIES'].median()
print(mednoi)
ldi['NO.OF_INQUIRIES'].replace(np.nan, mednoi, inplace=True)
1.0
In [121]:
#dropping useless columns
ldi=ldi.drop(['UniqueID','Current_pincode_ID','Date.of.Birth','Employee_code_ID'],axis=1)
ldi.head()
Out[121]:
disbursed_amount asset_cost ltv branch_id supplier_id manufacturer_id Age Employment.Type DisbursalDate State_ID ... SEC
0 505780 687630 89.55 67 22807 45 39271233 Salaried 2018-08-03 6 ...
1 532780 613600 88.39 67 22807 45 37624658 Self employed 2018-08-01 6 ...
2 523780 613600 88.39 67 22807 45 45336986 Self employed 2018-09-26 6 ...
3 463490 615000 76.42 67 22807 45 34852055 Salaried 2018-09-23 6 ...
4 435940 687630 79.77 67 22744 86 28731507 Self employed 2018-10-08 6 ...
5 rows x 38 columns
In [122]:
#transforming data types to their appropriate types
for x in obj_vars:
 ldi[x]=ldi[x].astype(str)
print(ldi.dtypes)
disbursed_amount float64
asset_cost float64
ltv float64
branch_id object
supplier_id object
manufacturer_id object
Age float64
Employment.Type object
DisbursalDate datetimestr
State_ID int64
MobileNo_Av1_Flag object
Aadhar_flag object
PAN_flag object
VoterID_flag object
Driving_flag object
Passport_flag object
Profile_strength object
PERFORM_CNS.SCORE int64
PERFORM_CNS.SCORE.DESCRPTION object
PRI.NO.OF.ACCTS int64
PRI.ACTIVE.ACCTS int64
PRI.OVERDUE.ACCTS int64
PRI.CURRENT.BALANCE float64
PRI.SANCTIONED.AMOUNT float64
PRI.DISBURSED.AMOUNT float64
SEC.NO.OF.ACCTS int64
SEC.ACTIVE.ACCTS int64
SEC.OVERDUE.ACCTS int64
SEC.CURRENT.BALANCE float64
SEC.SANCTIONED.AMOUNT float64
SEC.DISBURSED.AMOUNT float64
PRIMARY.INSTAL.AMT float64
SEC.INSTAL.AMT float64
NEW.ACCTS.IN.LAST.SIX.MONTHS int64
DELINQUENT.ACCTS.IN.LAST.SIX.MONTHS int64
Acc_age_months float64
NO.OF_INQUIRIES float64
loan_default object
dtype: object
In [123]:
#getting non categorical variables
non_cat=[]
for x in ldi.columns:
 if x in obj_vars:
 pass
 else:
 non_cat.append(x)
print(non_cat)
['disbursed_amount', 'asset_cost', 'ltv', 'Age', 'Employment.Type', 'DisbursalDate', 'Profile_strength', 'PERFORM_CNS.SCORE', 'PERFORM_CNS.SCORE.DESCRPTION', 'PRI.NO.OF.ACCTS', 'PRI.ACTIVE.ACCTS', 'PRI.OVERDUE.ACCTS', 'PRI.CURRENT.BALANCE', 'PRI.SANCTIONED.AMOUNT', 'PRI.DISBURSED.AMOUNT', 'SEC.NO.OF.ACCTS', 'SEC.ACTIVE.ACCTS', 'SEC.OVERDUE.ACCTS', 'SEC.CURRENT.BALANCE', 'SEC.SANCTIONED.AMOUNT', 'SEC.DISBURSED.AMOUNT', 'PRIMARY.INSTAL.AMT', 'SEC.INSTAL.AMT', 'NEW.ACCTS.IN.LAST.SIX.MONTHS', 'DELINQUENT.ACCTS.IN.LAST.SIX.MONTHS', 'Acc_age_months', 'NO.OF_INQUIRIES']
Step 3:Exploratory data analysis
In this step, we are going to explore the dataset. Perform hypothesis tests,bivariate analysis and check for correlation between variables.
In [124]:
#summary stats
ldi.describe().round(2)
Out[124]:
disbursed_amount asset_cost ltv Age PERFORM_CNS.SCORE PRI.NO.OF.ACCTS PRI.ACTIVE.ACCTS PRI.OVERDUE.ACCTS P
count 233154.00 233154.00 233154.00 233154.00 233154.00 233154.00 233154.00 233154.00
mean 53826.46 68860.58 79.88 39.01 289.46 2.44 1.04 0.16
std 6198.68 3150.76 5.07 9.83 338.37 5.22 1.94 0.55
min 39794.00 61330.00 68.88 22.46 0.00 0.00 0.00 0.00
25% 50078.00 67743.25 76.80 30.88 0.00 0.00 0.00 0.00
50% 53803.00 68763.00 79.77 37.27 0.00 0.00 0.00 0.00
75% 57363.00 69800.00 83.38 45.94 678.00 3.00 1.00 0.00
max 68862.00 76807.00 89.95 73.59 890.00 453.00 144.00 25.00
8 rows x 23 columns
In [125]:
#seeing correlation of numerical variables
corrdata =ldi[non_cat].corr()
corrdata.round(2)
Out[125]:
disbursed_amount asset_cost ltv Age PERFORM_CNS.SCORE PRI.NO.OF.ACCTS PRI.ACTIVE.ACCTS
disbursed_amount 1.00 0.27 0.30 -0.01 0.01 0.04 0.03
asset_cost 0.27 1.00 -0.04 -0.01 -0.03 -0.00 -0.01
ltv 0.30 -0.04 1.00 0.08 0.07 0.06 0.06
Age -0.01 -0.01 0.08 1.00 0.17 0.17 0.14
PERFORM_CNS.SCORE 0.01 -0.03 0.07 0.17 1.00 0.42 0.47
PRI.NO.OF.ACCTS 0.04 -0.00 0.06 0.17 0.42 1.00 0.75
PRI.ACTIVE.ACCTS 0.03 -0.01 0.06 0.14 0.47 0.75 1.00
PRI.OVERDUE.ACCTS 0.02 -0.00 0.03 0.14 0.11 0.35 0.38
PRI.CURRENT.BALANCE 0.00 -0.00 0.00 0.02 0.03 0.04 0.06
PRI.SANCTIONED.AMOUNT 0.00 -0.01 0.01 0.08 0.14 0.19 0.25
PRI.DISBURSED.AMOUNT 0.00 -0.00 0.02 0.06 0.09 0.17 0.22
SEC.NO.OF.ACCTS -0.01 -0.02 0.00 0.02 0.06 0.06 0.06
SEC.ACTIVE.ACCTS -0.01 -0.02 0.00 0.02 0.05 0.05 0.06
SEC.OVERDUE.ACCTS -0.01 -0.01 0.00 0.02 0.04 0.03 0.04
SEC.CURRENT.BALANCE -0.00 0.00 -0.00 0.00 0.00 0.01 0.01
SEC.SANCTIONED.AMOUNT 0.00 0.00 -0.00 0.00 -0.00 0.01 0.01
SEC.DISBURSED.AMOUNT 0.00 0.00 0.00 0.00 -0.00 0.00 0.00
PRIMARY.INSTAL.AMT 0.01 -0.00 0.02 0.08 0.15 0.20 0.23
SEC.INSTAL.AMT -0.00 -0.00 0.00 0.00 0.02 0.01 0.02
NEW.ACCTS.IN.LAST.SIX.MONTHS 0.03 -0.01 0.06 0.03 0.35 0.54 0.70
DELINQUENT.ACCTS.IN.LAST.SIX.MONTHS 0.01 -0.00 0.02 0.08 0.16 0.31 0.38
Acc_age_months -0.01 -0.00 -0.01 0.10 0.10 0.01 0.06
NO.OF_INQUIRIES 0.02 -0.01 0.04 0.00 0.10 0.15 0.16
23 rows x 23 columns
Bivariate Analysis
We are going to be analysing two variables with respect to one another here.
In [126]:
#are high enquires a red flag
plt.figure(figsize=(12,6))
ldi.boxplot(by="PERFORM_CNS.SCORE.DESCRPTION",column="NO.OF_INQUIRIES", grid = False)
<AxesSubplot: title='(center: 'NO.OF_INQUIRIES', xlabel='PERFORM_CNS.SCORE.DESCRPTION'>
<Figure size 120x600 with 0 Axes>
Boxplot grouped by PERFORM_CNS.SCORE.DESCRPTION
NO.OF_INQUIRIES

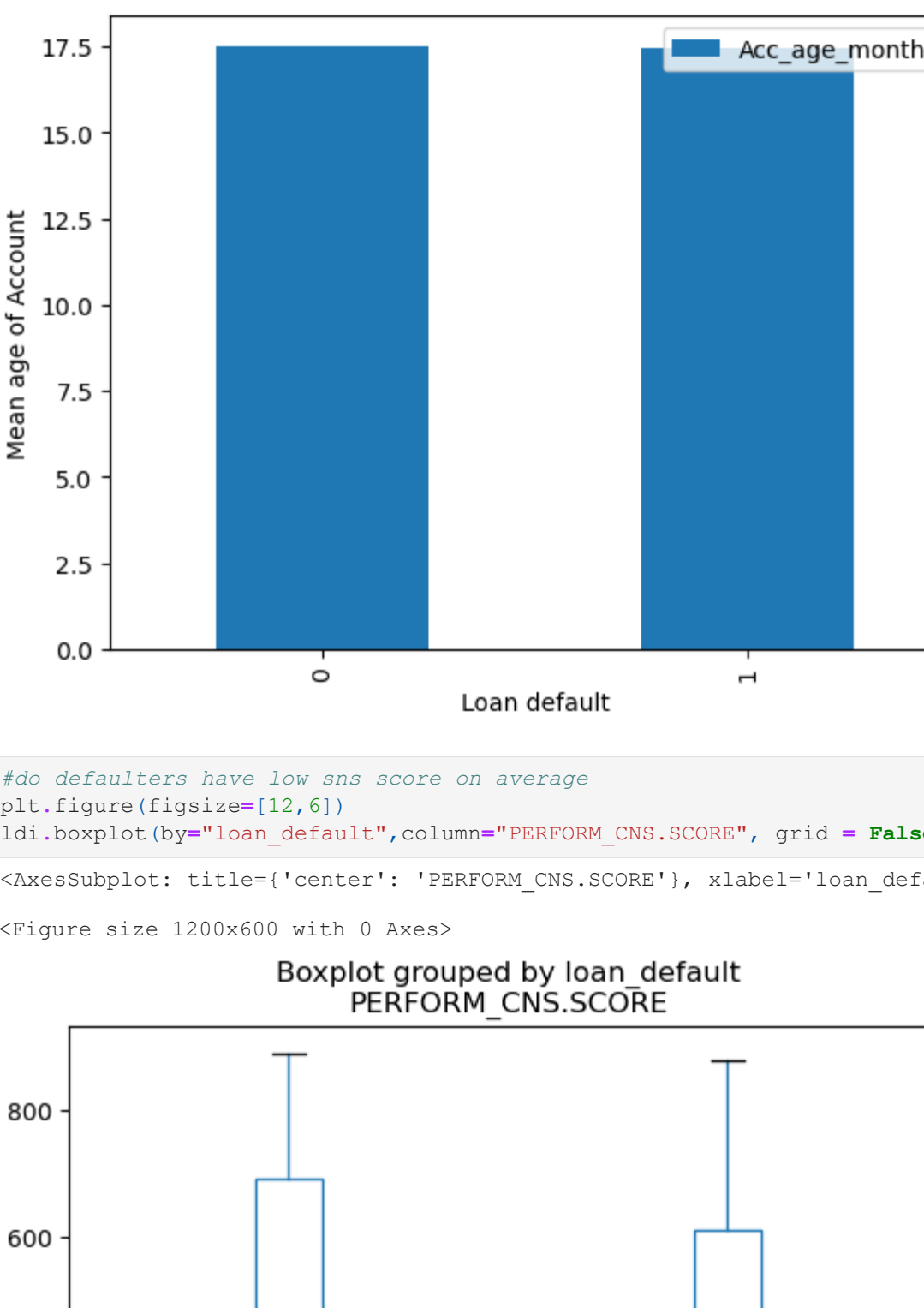
In [127]:
#checking if enquiries differs according to the risk category'
anov=ldi['NO.OF_INQUIRIES']
riskcat=ldi['PERFORM_CNS.SCORE.DESCRPTION']
mod = ols("enqnt ~ riskcat", data=ldi).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
Out[127]:
df sum_sq mean_sq F PR(>F)
riskcat 6.0 380.521844 63.420307 622.73396 0.0
Residual 233147.0 23744.095126 0.101842 NaN NaN
In [128]:
#seeing the number of enquiries by credit risk'
grouped = ldi.groupby('PERFORM_CNS.SCORE.DESCRPTION').agg({'NO.OF_INQUIRIES': 'mean'})
grouped.plot(kind='bar')
plt.xlabel('Credit risk type')
plt.ylabel('Mean number of enquiries')
plt.show()

In [129]:
#discrepancy calculation
ldi['pri_diff']=ldi['PRI.SANCTIONED.AMOUNT']-ldi['PRI.DISBURSED.AMOUNT']
ldi['sec_diff']=ldi['SEC.SANCTIONED.AMOUNT']-ldi['SEC.DISBURSED.AMOUNT']
ldi.head()
Out[129]:
disbursed_amount asset_cost ltv branch_id supplier_id manufacturer_id Age Employment.Type DisbursalDate State_ID ... SEC
0 505780 687630 89.55 67 22807 45 39271233 Salaried 2018-08-03 6 ...
1 532780 613600 88.39 67 22807 45 37624658 Self employed 2018-08-01 6 ...
2 523780 613600 88.39 67 22807 45 45336986 Self employed 2018-09-26 6 ...
3 463490 615000 76.42 67 22807 45 34852055 Salaried 2018-09-23 6 ...
4 435940 687630 79.77 67 22744 86 28731507 Self employed 2018-10-08 6 ...
5 rows x 40 columns
In [130]:
#difference btw sanctioned amount disbursed amount btw primary and secondary across loan default types
ldi['pri_diff']=ldi['PRI.SANCTIONED.AMOUNT']-ldi['PRI.DISBURSED.AMOUNT']
ldi['sec_diff']=ldi['SEC.SANCTIONED.AMOUNT']-ldi['SEC.DISBURSED.AMOUNT']
ldi.head()
value1=ldi['pri_diff'].mean()
value2=ldi['sec_diff'].mean()
size of width of the bars
bar_width = 0.35
set the positions of the bars on the x-axis
r1 = np.arange(len(ldi))
r2 = [x + bar_width for x in r1]
Create the bar chart
plt.bar(r1, value1, width=bar_width, edgecolor='white', label='primary amount difference')
plt.bar(r2, value2, width=bar_width, edgecolor='white', label='secondary amount difference')
Add ticks on the middle of the group bars
plt.xticks([len(ldi)/2], label='value')
Add legend
plt.legend()
Show the chart
plt.show()

In [131]:
#which is the most used mode for registration
anov=ldi['MobileNo_Av1_Flag','Aadhar_flag','PAN_flag','VoterID_flag','Driving_flag','Passport_flag']
modes.head()
counts = modes.apply(lambda x: x.value_counts()[1])
Create the bar chart
plt.bar(counts.index, counts.values)
plt.xlabel('Mode of identity')
plt.ylabel('Count')
plt.title('Whats the most used mode by customers')
Show the chart
plt.show()

In [132]:
#is account age responsible for default
plt.figure(figsize=(12,6))
ldi.boxplot(by="loan_default",column="Acc_age_months", grid = False)
<AxesSubplot: title='(center: 'Acc_age_months', xlabel='loan_default'>
<Figure size 120x600 with 0 Axes>
Boxplot grouped by loan_default
Acc_age_months

In [133]:
#checking if account age differs according to the loan default
anov=ldi['Acc_age_months']
ldi['loan_default']
mod = ols("age ~ ld", data=ldi).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
Out[133]:
df sum_sq mean_sq F PR(>F)
ld 1.0 9.137824e+01 91378235 7.822513 0.00516
Residual 233152.0 2.723525e+06 11.681442 NaN NaN
In [134]:
#seeing the naverage age of bank accounts according to loan default
grouped = ldi.groupby('loan_default').agg({'Acc_age_months': 'mean'})
grouped.plot(kind='bar')
plt.xlabel('loan default')
plt.ylabel('Mean age of account')
plt.show()



```
In [135]: #do defaulters have low sns score on average
plt.figure(figsize=(12,6))
l=ldi.boxplot(by="loan_default",column="PERFORM_CNS.SCORE", grid = False)
```

```
Out[135]: <AxesSubplot: title='center': 'PERFORM_CNS.SCORE', xlabel='loan_default'>
```

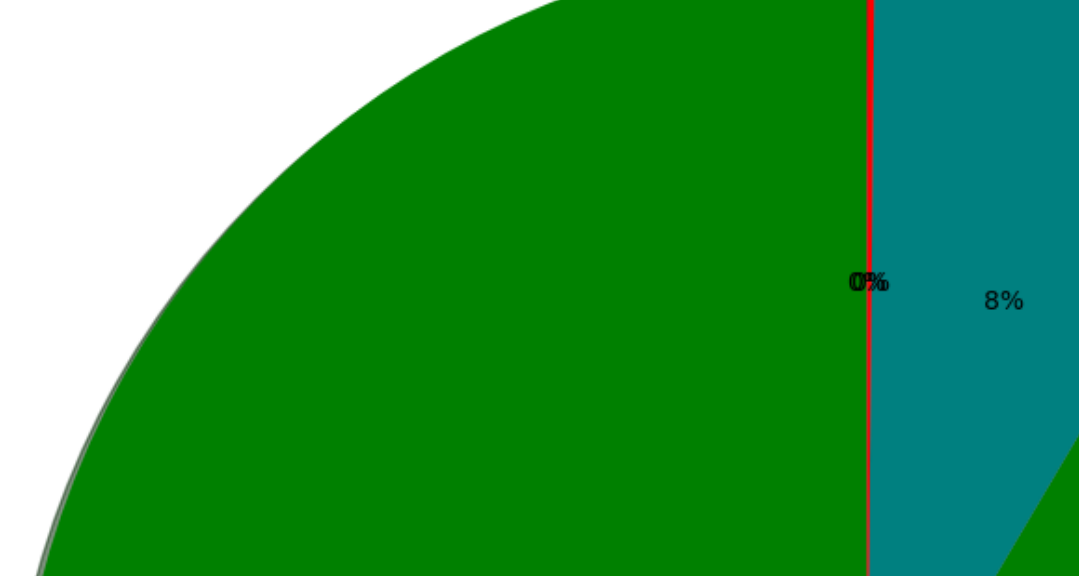
```
<Figure size 1200x600 with 0 Axes>
```



```
In [136]: #checking if score differs according to the loan default
score=ldi['PERFORM_CNS.SCORE']
l=ldi['loan_default']
mod = ols("score ~ ldi",data=ldi).fit()
anova_table=sns.stats.anova_lm(mod)
anova_table
```

	df	sum_sq	mean_sq	F	PR(>F)
ldi	1	0.8958388e+07	8.958388e+07	785.040202	1.869421e-172
Residual	2331520	2.660585e+10	1.141138e+05	NaN	NaN

```
In [137]: #seeing the naverage age of bank accounts according to loan default
grouped = ldi.groupby('loan_default').agg({'PERFORM_CNS.SCORE': 'mean'})
```

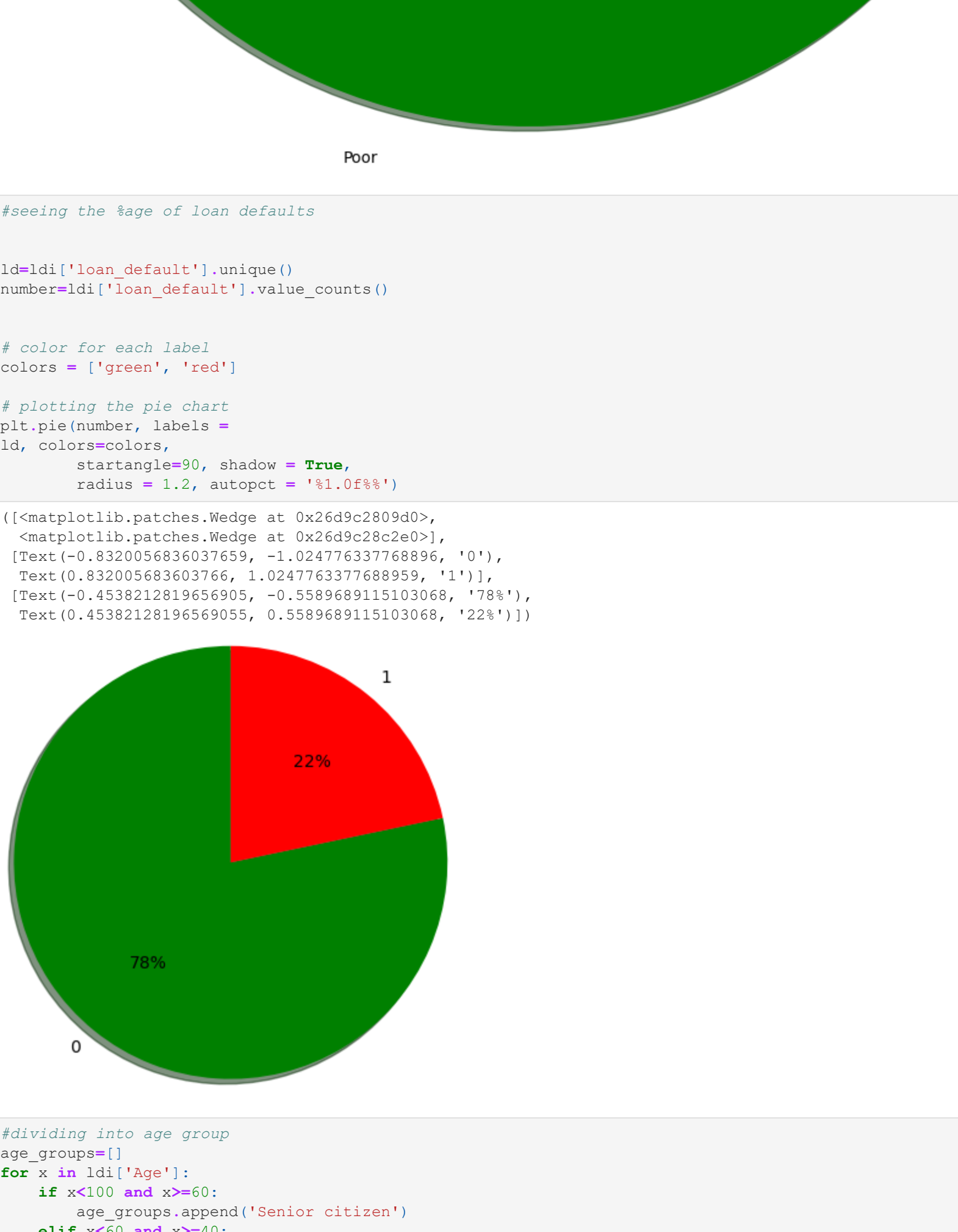


```
In [138]: #What is the profile strength of majority of people
```

```
p=ldi['Profile_strength'].unique()
numbers=ldi['Profile_strength'].value_counts()

# color for each label
colors = ['green','teal','red','blue']

# plotting the pie chart
plt.pie(number, labels =
p, colors=colors,
startangle=90, shadow = True,
radius = 3.1, autopct = '%1.0f%%')
```

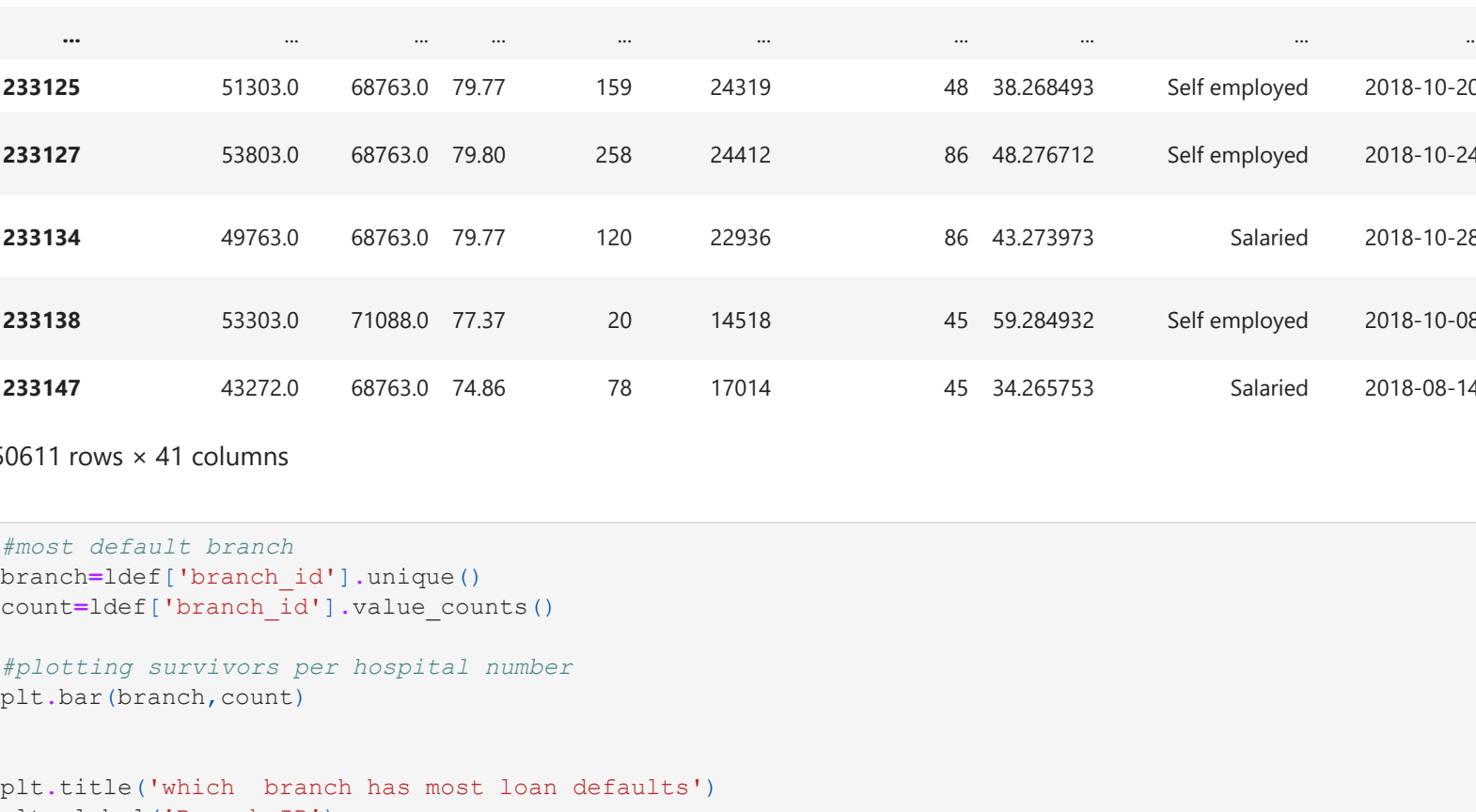


```
In [139]: #seeing the sage of loan defaults
```

```
l=ldi['loan_default'].unique()
numbers=ldi['loan_default'].value_counts()

# color for each label
colors = ['green', 'red']

# plotting the pie chart
plt.pie(number, labels =
ld, colors=colors,
startangle=90, shadow = True,
radius = 1.2, autopct = '%1.0f%%')
```



```
In [140]: #dividing into age group
age_group=[]
for x in range(1,age):
    if x<100 and x>=60:
        age_group.append('Senior citizen')
    elif x<60 and x>=40:
        age_group.append('Middle Aged')
    elif x<40 and x>=25:
        age_group.append('Adult')
    else:
        age_group.append('Youth')
```

```
ldi['age_group']=age_group
```

```
In [141]: #taking subset of default
lde=ldi[ldi.loan_default==1]
lde
```

```
Out[141]:
```

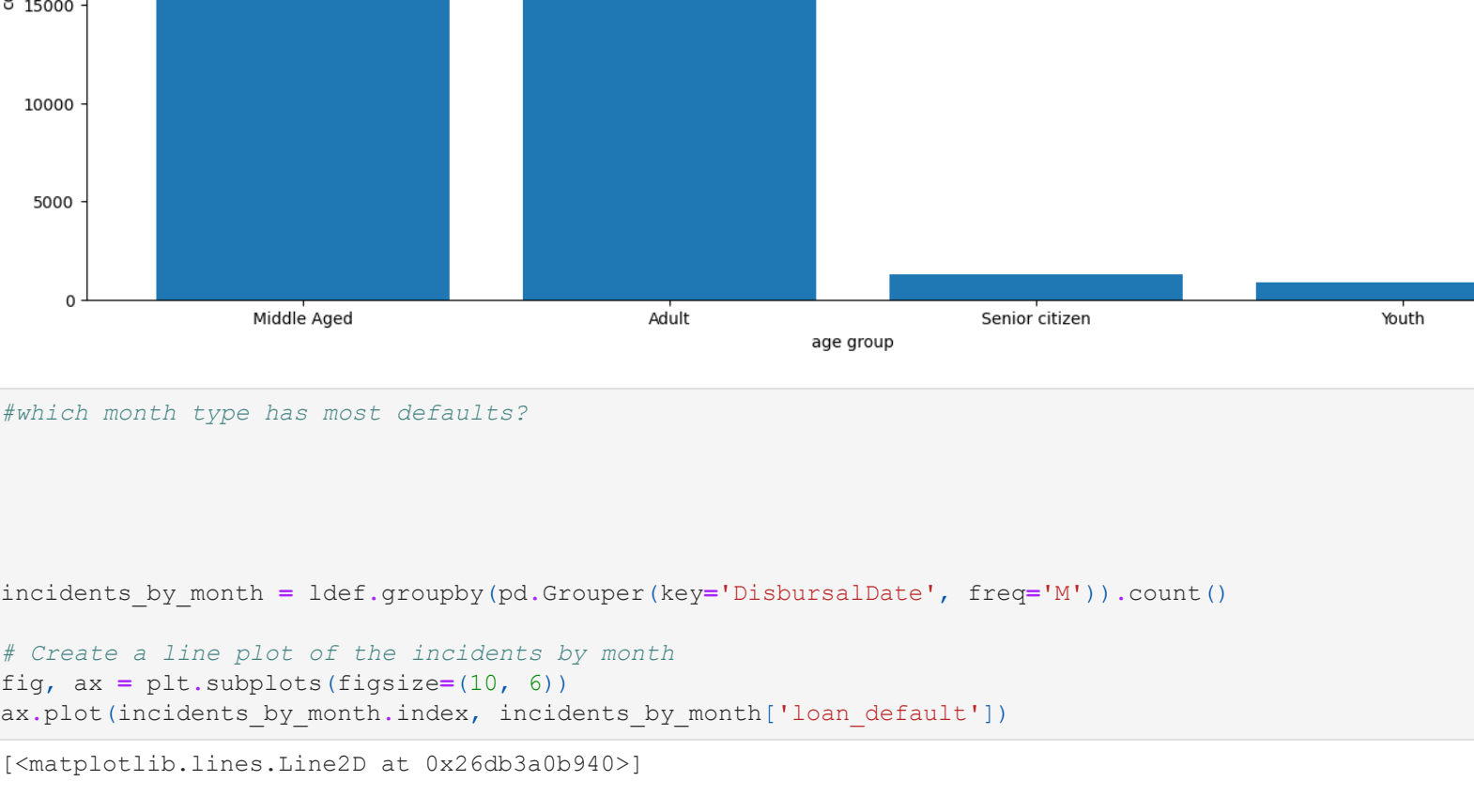
	disbursed_amount	asset_cost	ltv	branch_id	supplier_id	manufacturer_id	Age	Employment.Type	DisbursalDate	State_ID	...
2	537130	687630	88.39	67	22807	45	45.336986	Self employed	2018-09-26	6	...
7	577130	617800	89.83	67	22807	45	54.695890	Self employed	2018-08-16	6	...
8	577130	680400	86.27	67	22807	45	47.276712	Self employed	2018-10-10	6	...
12	505780	687630	89.55	67	22807	45	26.638356	Salaried	2018-08-06	6	...
18	555130	679500	83.15	67	22807	45	28.841096	Salaried	2018-10-04	6	...
...
233125	513030	687630	79.77	159	24319	48	38.269493	Self employed	2018-10-20	7	...
233127	538030	687630	79.80	258	24412	86	48.276712	Self employed	2018-10-24	16	...
233134	497630	687630	79.37	120	22936	86	43.273973	Salaried	2018-10-28	12	...
233138	532030	710880	79.77	20	14518	45	59.284932	Self employed	2018-10-08	5	...
233147	439540	687630	74.86	78	17014	45	34.265753	Salaried	2018-08-14	4	...

50611 rows x 41 columns

```
In [142]: #most default branch
branch=ldi['branch_id'].unique()
count=ldi['branch_id'].value_counts()

#plotting survivors per hospital number
plt.bar(branch,count)
```

```
plt.title('which branch has most loan defaults')
plt.xlabel('branch ID')
plt.ylabel('count')
plt.rcParams['figure.figsize'] = (17, 7)
plt.show()
```



```
In [143]: #most default state
state=ldi['State_ID'].unique()
count=ldi['State_ID'].value_counts()

#plotting survivors per hospital number
plt.bar(state,count)
```

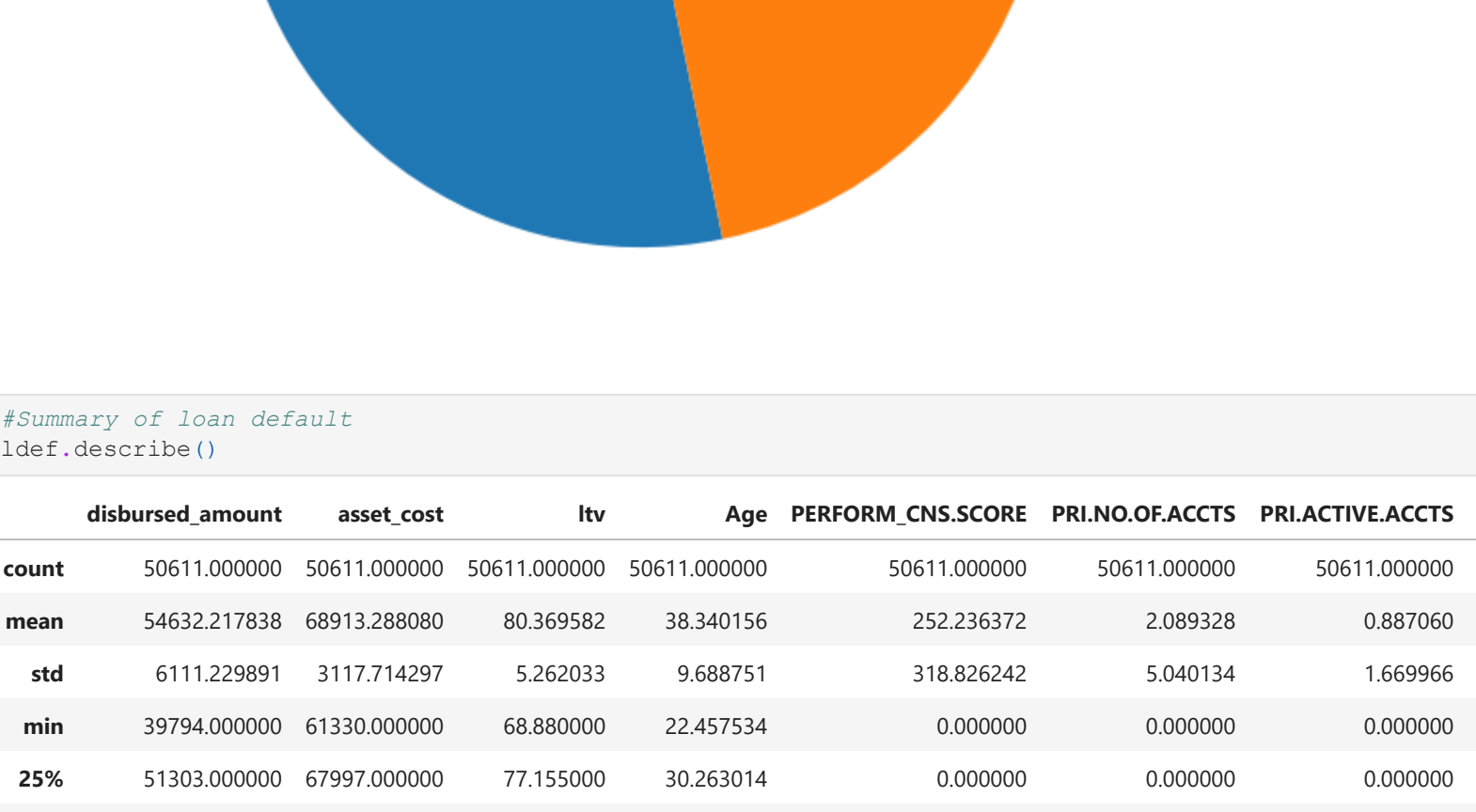
```
plt.title('which State has most loan defaults')
plt.xlabel('state')
plt.ylabel('count')
plt.rcParams['figure.figsize'] = (17, 7)
plt.show()
```



```
In [144]: #which age group is likely to default on loans?
age_group=ldi['age_group'].unique()
count=ldi['age_group'].value_counts()

#plotting survivors per hospital number
plt.bar(age_group,count)
```

```
plt.title('which age group is likely to default on loans?')
plt.xlabel('age group')
plt.ylabel('count')
plt.rcParams['figure.figsize'] = (17, 7)
plt.show()
```

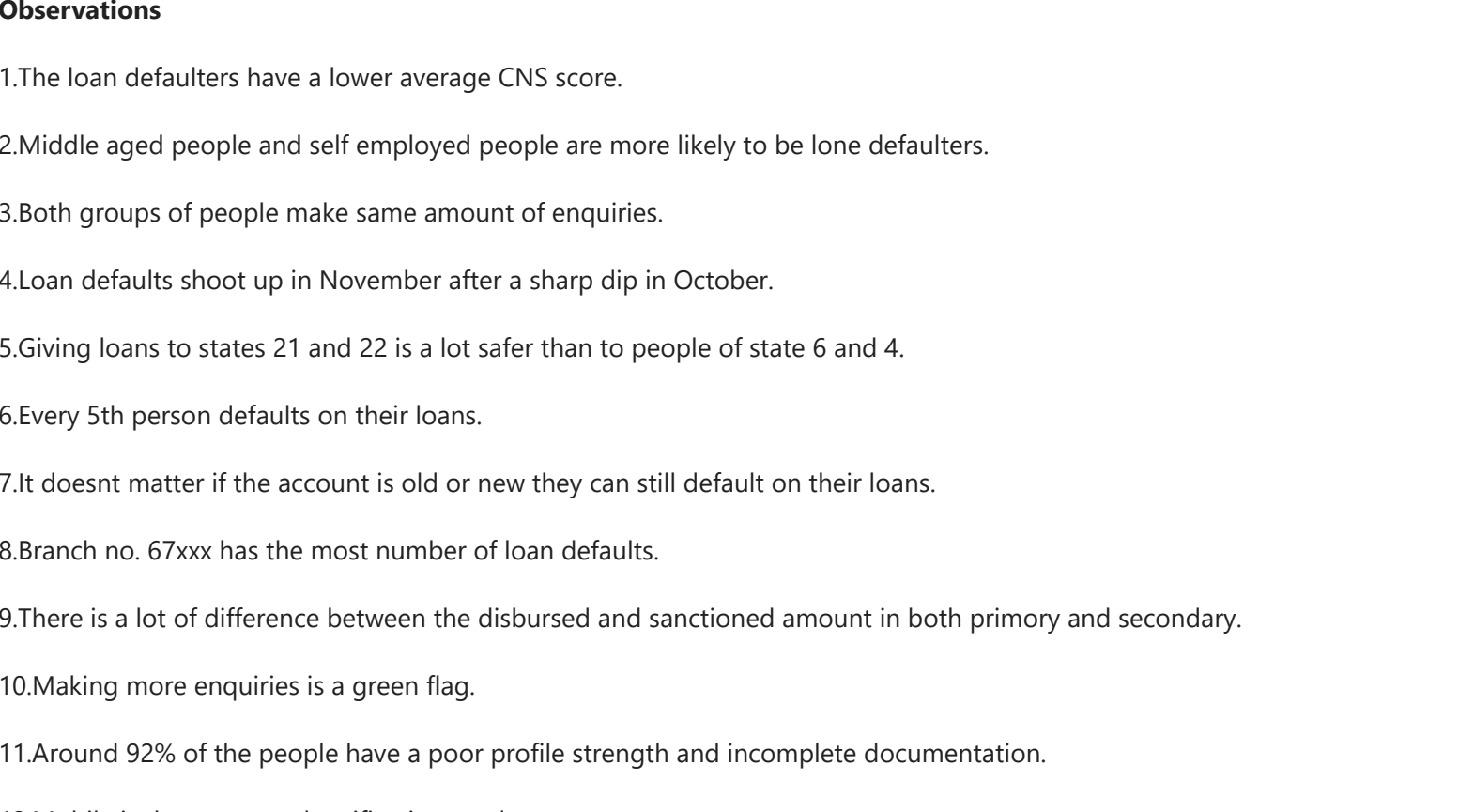


```
In [145]: #which month type has most defaults?

incidents_by_month = lde.groupby(pd.Grouper(key='DisbursalDate', freq='M')).count()

# Create a line plot of the incidents by month
fig = plt.subplots(figsize=(10, 6))
ax=plt.plot(incidents_by_month.index, incidents_by_month['loan_default'])

[matplotlib.lines.Line2D at 0x26db3a0b940]
```



```
In [146]: #What employment type are majorly in loan default?
emp_type=ldi['Employment.Type'].unique()
numbers=ldi['Employment.Type'].value_counts()

#employment type percentage
plt.pie(numbers, labels = emp_type, startangle = 75, autopct='%1.2f%%')
plt.show()
```



```
In [147]: #Summary of loan default
lde.describe()
```

	disbursed_amount	asset_cost	ltv	Age	PERFORM_CNS.SCORE	PRI.NO.OF.ACCTS	PRI.ACTIVE.ACCTS	PRI.OVERDUE.ACCTS	PRI.OVERDUE.DUE
count	50611.000000	50611.000000	50611.000000	50611.000000	50611.000000	50611.000000	50611.000000	50611.000000	50611.000000
mean	54632.217838	68913.288080	80.369582	38.340156	252.236372	2.089328	0.887060	0.887060	0.887060
std	6111.229891	3117.714297	5.262053	9.688751	318.826242	5.040134	1.669966	1.669966	1.669966
min	39794.000000	61330.000000	68.880000	22.457534	0.000000	0.000000	0.000000	0.000000	0.000000
25%	51303.000000	67997.000000	77.155000	30.263014	0.000000	0.000000	0.000000	0.000000	0.000000
50%	53803.000000	68763.000000	79.770000	36.268493	0.000000	0.000000	0.000000	0.000000	0.000000
75%	5813.000000	68600.000000	84.540000	45.056164	610.000000	2.000000	1.000000	1.000000	1.000000
max	68882.000000	76807.000000	89.950000	68.758904	879.000000	453.000000	35.000000	35.000000	35.000000

8 rows x 25 columns

Observations

- 1.The loan defaulters have a lower average CNS score.
- 2.Middle aged people and self employed people are more likely to be lone defaulters.
- 3.Both groups of people make same amount of enquiries.
- 4.Loan defaults shoot up in November after a sharp dip in October.
- 5.Giving loans to states 21 and 22 is a lot safer than to people of state 6 and 4.
- 6.Every 5th person defaults on their loans.
- 7.It doesn't matter if the account is old or new they can still default on their loans.
- 8.Branch no. 67xxx has the most number of loan defaults.
- 9.There is a lot of difference between the disbursed and sanctioned amount in both primary and secondary.
- 10.Making more enquiries is a green flag.
- 11.Around 92% of the people have a poor profile strength and incomplete documentation.
- 12.Mobile is the most used verification mode.
- 13.The amount disbursed has a moderate positive correlation with ltv ratio and asset value.

Step 4: Model Building

Now we are going to build the Decision Tree classifier to predict who is more likely to default on a loan.

```
In [148]: #creating target variable which is the outcome of the loan default
target=ldi['loan_default']
```

```
In [149]: #removing target variable from main dataframe
ldi= ldi.drop(['loan_default'], axis=1)
```

```
In [150]: #flattening categorical variables
obj_vars=['Employment.Type','age_group','branch_id','supplier_id','manufacturer_id','PERFORM_CNS.SCORE.DESCRIP
'MobileNo_Avi_Flag','Aadhar_Flag','PAN_Flag','VoterID_Flag','Driving_Flag','Passport_Flag','Profile_s
```

```
In [151]: #fitting the model
ldi[obj_vars]=['Employment.Type','age_group','branch_id','supplier_id','manufacturer_id','PERFORM_CNS.SCORE.DES
'MobileNo_Avi_Flag','Aadhar_Flag','PAN_Flag','VoterID_Flag','Driving_Flag','Passport_Flag','Profile_s
```

```
In [152]: #creating dummy variables

for category in ldi[obj_vars]:
    ldi[category] = pd.get_dummies(ldi[category])
```

```
In [153]: #checking missing values
ldi.isnull().sum()
```

	disbursed_amount	asset_cost	ltv	branch_id	supplier_id	manufacturer_id	Age	Employment.Type	DisbursalDate	State_ID	...	SEC
0	505780	687630	89.55	1	1	1	39.271233	1	2018-08-03	1	...	1
1	532780	613600	89.63	1	1	1	37.624658	1	2018-08-01	1	...	1
2	523780	687630	88.39	1	1	1	45.336986	1	2018-09-26	1	...	1
3	463490	615000	76.42	1	1	1	34.852055	1	2018-09-23	1	...	1
4	435940	687630	79.77	1	1	1	28.731507	1	2018-10-08	1	...	1

5 rows x 39 columns

```
In [155]: #checking shape of transformed dataset
ldi.shape
```

```
Out[155]: (233154, 40)
```

```
In [156]: #getting classifier
dct=DecisionTreeClassifier()
```

```
Out[156]:
```

```
In [157]: #initializing label encoder
le=LabelEncoder()
```

```
Out[157]:
```

```
In [158]: #standard scaler initiated
sc=StandardScaler()
```

```
In [159]: #transforming target variable
target=ldi.fit_transform(target)
```

```
Out[159]: array([0, 0, 1, ..., 0, 0, 0])
```

```
In [160]: #dropping useless columns again
ldi=ldi.drop(['DisbursalDate'],axis=1)
ldi.head()
```

	disbursed_amount	asset_cost	ltv	branch_id	supplier_id	manufacturer_id	Age	Employment.Type	State_ID	MobileNo_Avi_Flag	...
0	505780	687630	89.55	1	1	1	39.271233	1	1	1	...
1	532780	613600	89.63	1	1	1	37.624658	1	1	1	...
2	523780	687630	88.39	1	1	1	45.336986	1	1	1	...
3	463490	615000	76.42	1	1	1	34.852055	1	1	1	...
4	435940	687630	79.77	1	1	1	28.731507	1	1	1	...

5 rows x 39 columns

```
In [161]: #converting to dataframe
ldi= pd.DataFrame(ldi)
target= pd.DataFrame(target)
```

```
In [162]: #strating X and y
X = ldi.values
y = target.values
```

```
In [163]: #train test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size =0.25, random_state=1)
```

```
In [164]: #using standard scaler
X_train=sc.fit_transform(X_train)
X_test=sc.transform(X_test)
```

```
In [165]: #fitting the model
dct.fit(X_train,y_train)
```

```
Out[165]:
```

```
In [166]: #predict
ypred= dct.predict(X_test)
```

```
Out[166]: array([0, 0, 0, ..., 0, 1, 0])
```

```
In [167]: #accuracy
print("Accuracy of the model: {}%".format(accuracy_score(y_test,ypred)*100))
```

Accuracy of the model: 66.659232452298%

```
In [168]: #confusion matrix
cm=confusion_matrix(y_test,ypred)
print(cm)
```

```
[[35597  9974]
 [ 9460 32589]]
```

```
In [169]: #printing metrics
print(classification_report(y_test,ypred))
```

	precision	recall	f1-score	support
0	0.79	0.78	0.79	45571
1	0.25	0.26	0.25	12718
accuracy			0.67	58289
macro avg	0.52	0.52	0.52	58289
weighted avg	0.67	0.67	0.69	58289

```
In [170]: #using Random Forest
rf=RandomForestClassifier()
```

```
In [171]: #fitting the model
rf.fit(X_train,y_train)
```

```
Out[171]:
```

```
In [172]: #predict
ypred= rf.predict(X_test)
```

```
Out[172]: array([0, 0, 0, ..., 0, 0, 0])
```

```
In [173]: #accuracy
print("Accuracy of the model: {}%".format(accuracy_score(y_test,ypred)*100))
```

Accuracy of the model: 78.1605487124501%

```
In [174]: #confusion matrix
cm=confusion_matrix(y_test,ypred)
print(cm)
```

```
[[43993  1578]
 [12096  6229]]
```

```
In [175]: #printing metrics
print(classification_report(y_test,ypred))
```

	precision	recall	f1-score	support
0	0.78	1.00	0.88	45571
1	0.43	0.00	0.01	12718
accuracy			0.78	58289
macro avg	0.61	0.50	0.44	58289
weighted avg	0.71	0.78	0.69	58289

Conclusion

We have used three models to classify the people into loan defaulter or not. The third model has been found to be the most accurate(almost 80%) and hence this is the best model to classify people into defaulters and non defaulters.