

Exploratory Data Analysis

This notebook reviews the overall dataset, and gives a breakdown of every metric both as a stand-alone metric and also separated by the deposit variable.

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: pd.set_option('display.max_columns',None)

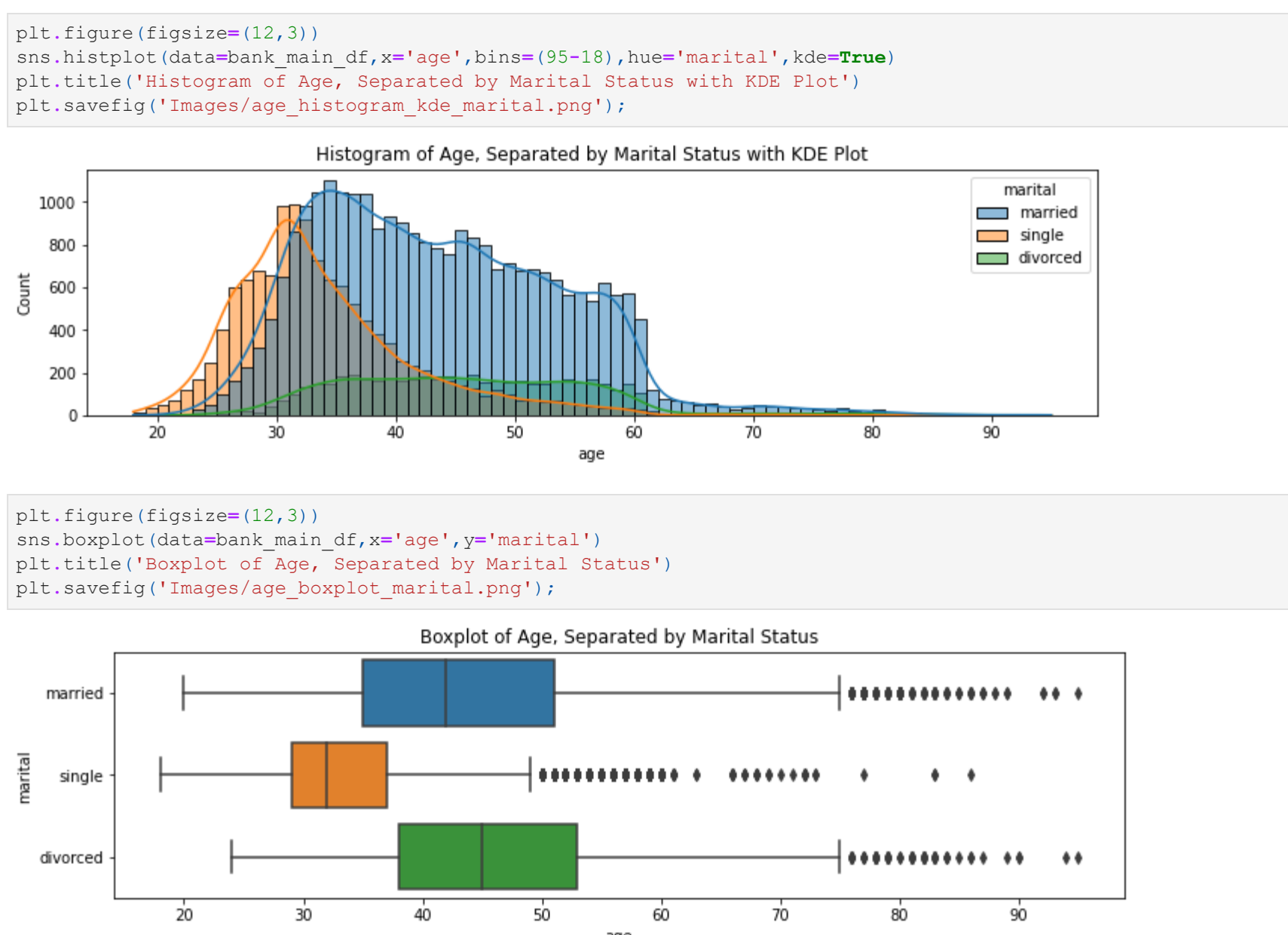
In [3]: bank_main_df = pd.read_csv('./Dataset_1_bank Marketing/bank_marketing.csv',delimiter=',')
bank_main_df.head(5)

Out[3]:
   age  job  marital  education  default  balance  housing  loan  contact  day  month  duration  campaign  pdays  previous  poutc
0  58.0  management  married  tertiary  no  2143  yes  no  unknown  5  may  261  1  -1  0  no
1  44.0  technician  single  secondary  no  29  yes  no  unknown  5  may  151  1  -1  0  no
2  33.0  entrepreneur  married  secondary  no  2  yes  yes  no  unknown  5  may  76  1  -1  0  no
3  47.0  blue-collar  married  secondary  no  1506  yes  no  unknown  5  may  92  1  -1  0  no
4  33.0  unknown  single  unknown  no  1  no  no  NaN  5  may  198  1  -1  0  no

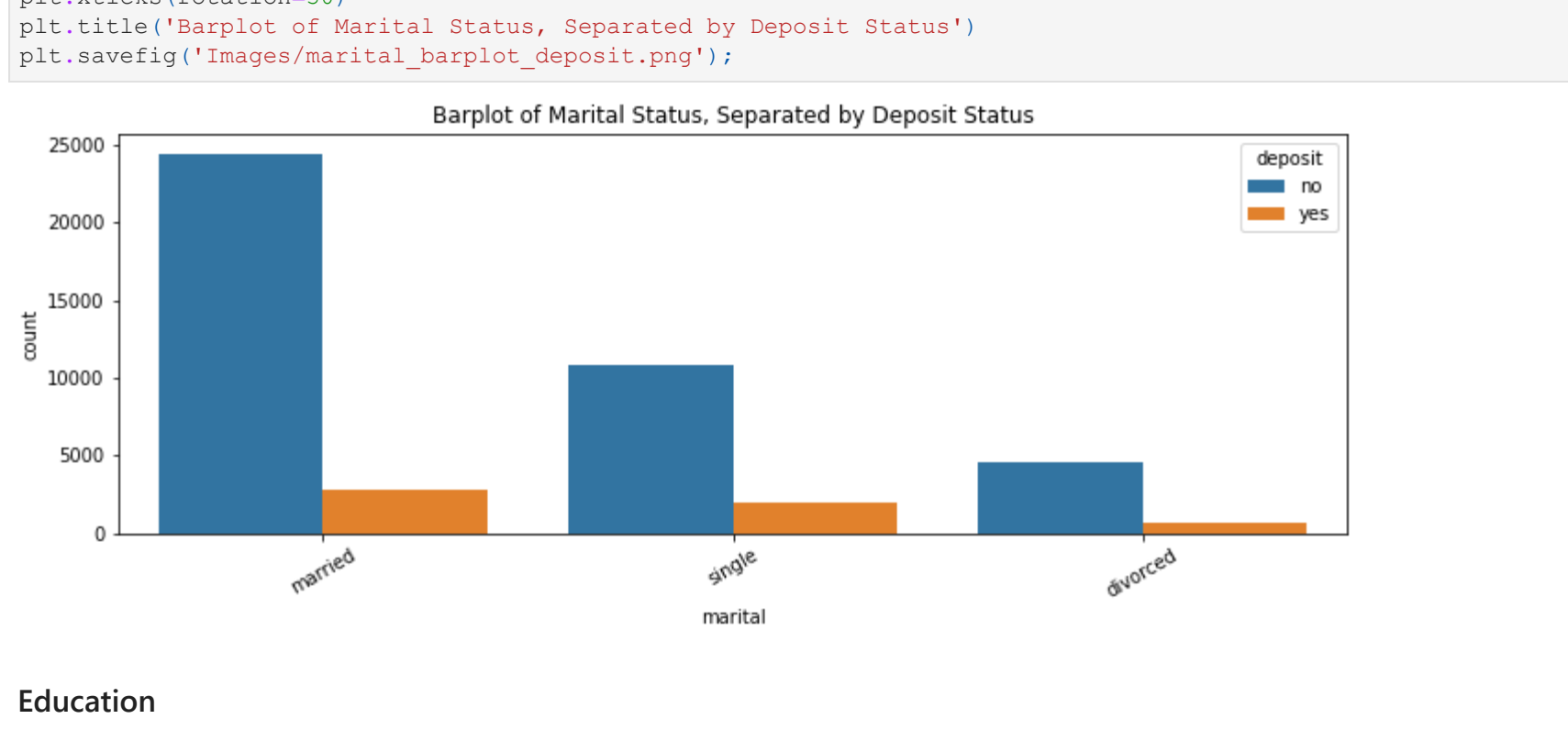
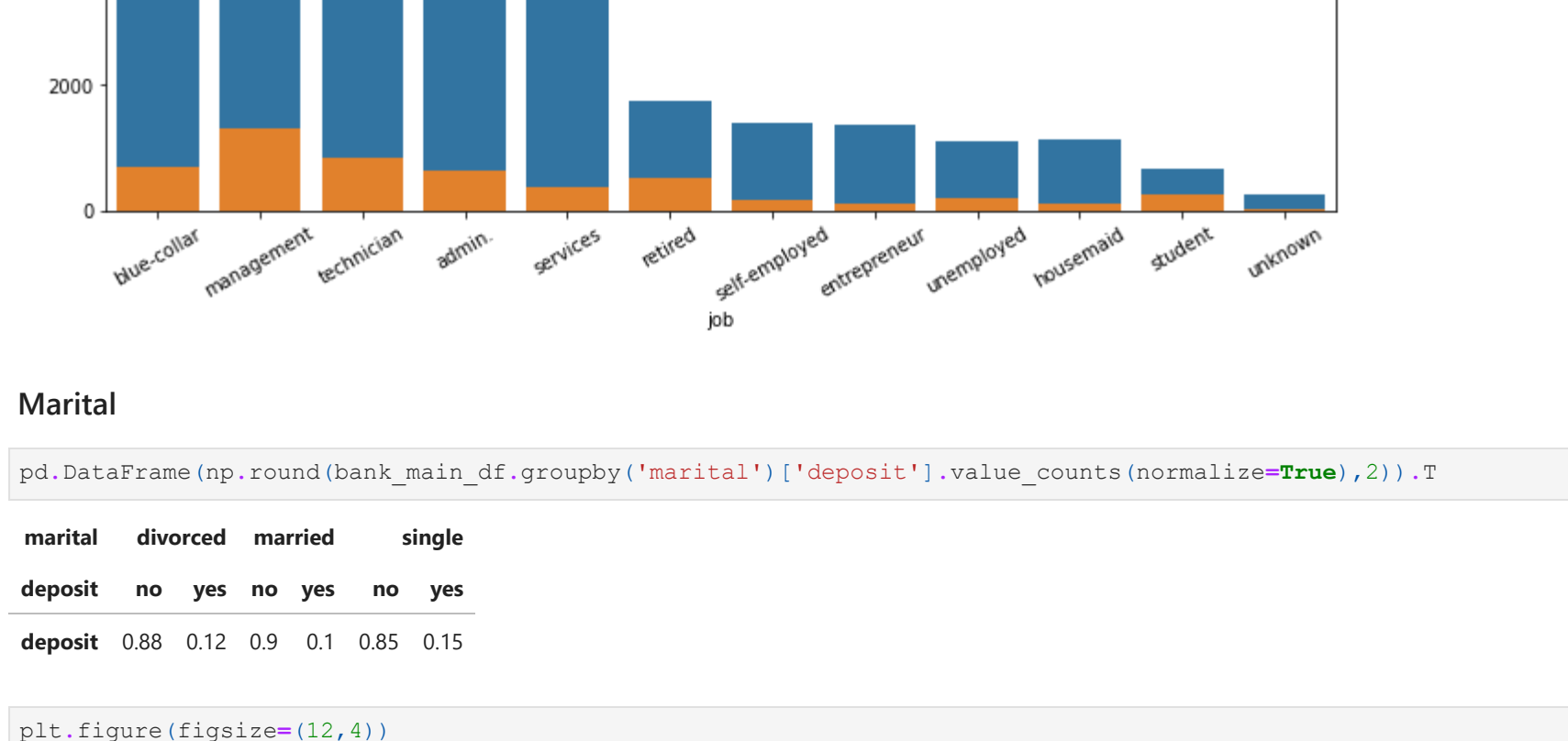
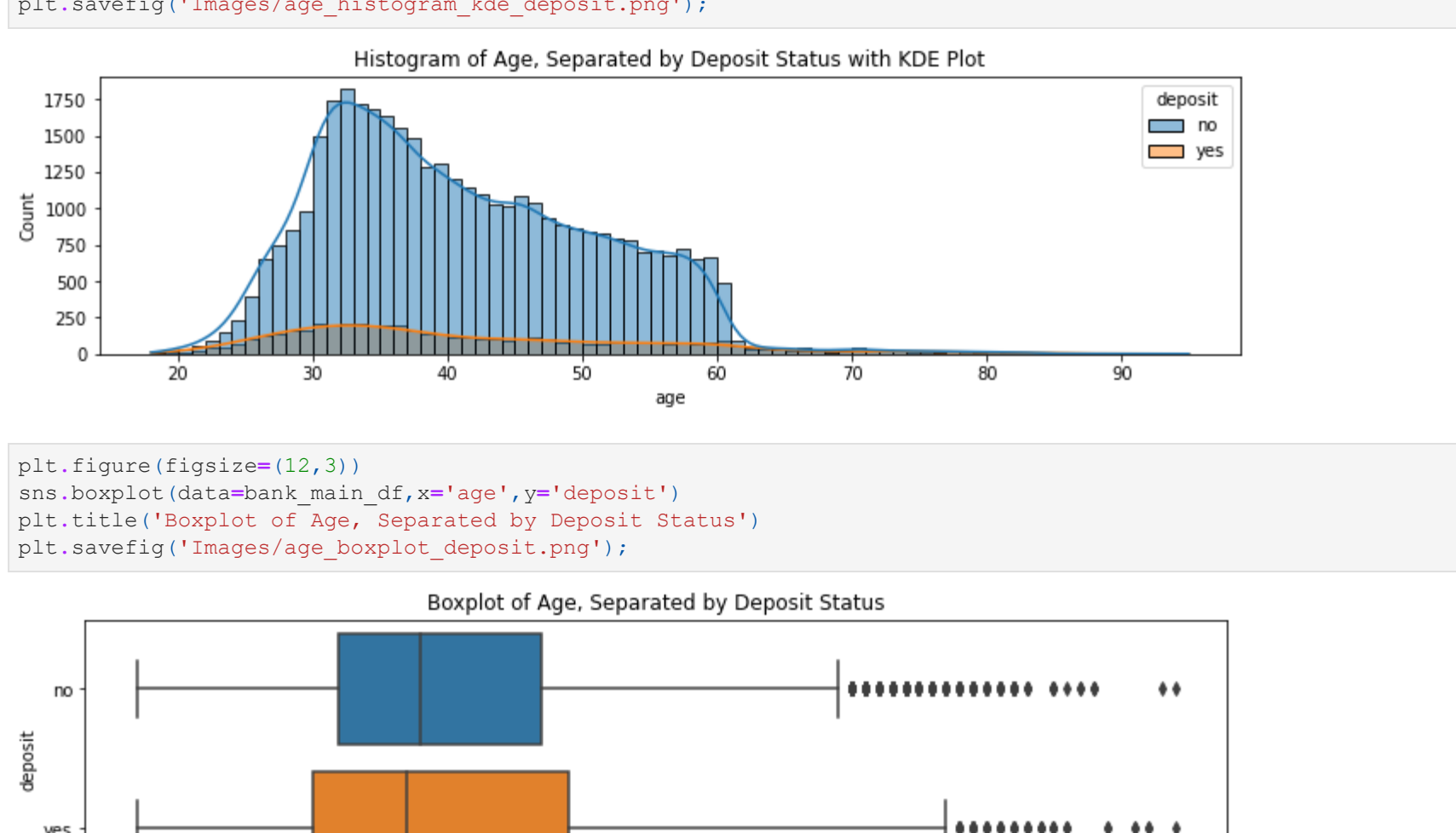
In [4]: bank_main_df.describe()

Out[4]:
      age      balance      day      duration      campaign      pdays      previous
count  43872.000000  45211.000000  45211.000000  45211.000000  45211.000000  45211.000000  45211.000000
mean    40.924781    1362.272058    15.806419    258.163080    2.763841    40.197828    0.580323
std     10.610835     3044.765829     8.322476    257.527812     3.098021    100.128746    2.303441
min     18.000000    -8019.000000     1.000000     0.000000     1.000000    -1.000000    0.000000
25%     33.000000     72.000000     8.000000    103.000000    1.000000    -1.000000    0.000000
50%     46.000000    148.000000    16.000000    180.000000    2.000000    -1.000000    0.000000
75%     58.000000    1028.000000    21.000000    319.000000    3.000000    -1.000000    0.000000
max     95.000000   142127.000000    31.000000    4918.000000    63.000000    871.000000    275.000000

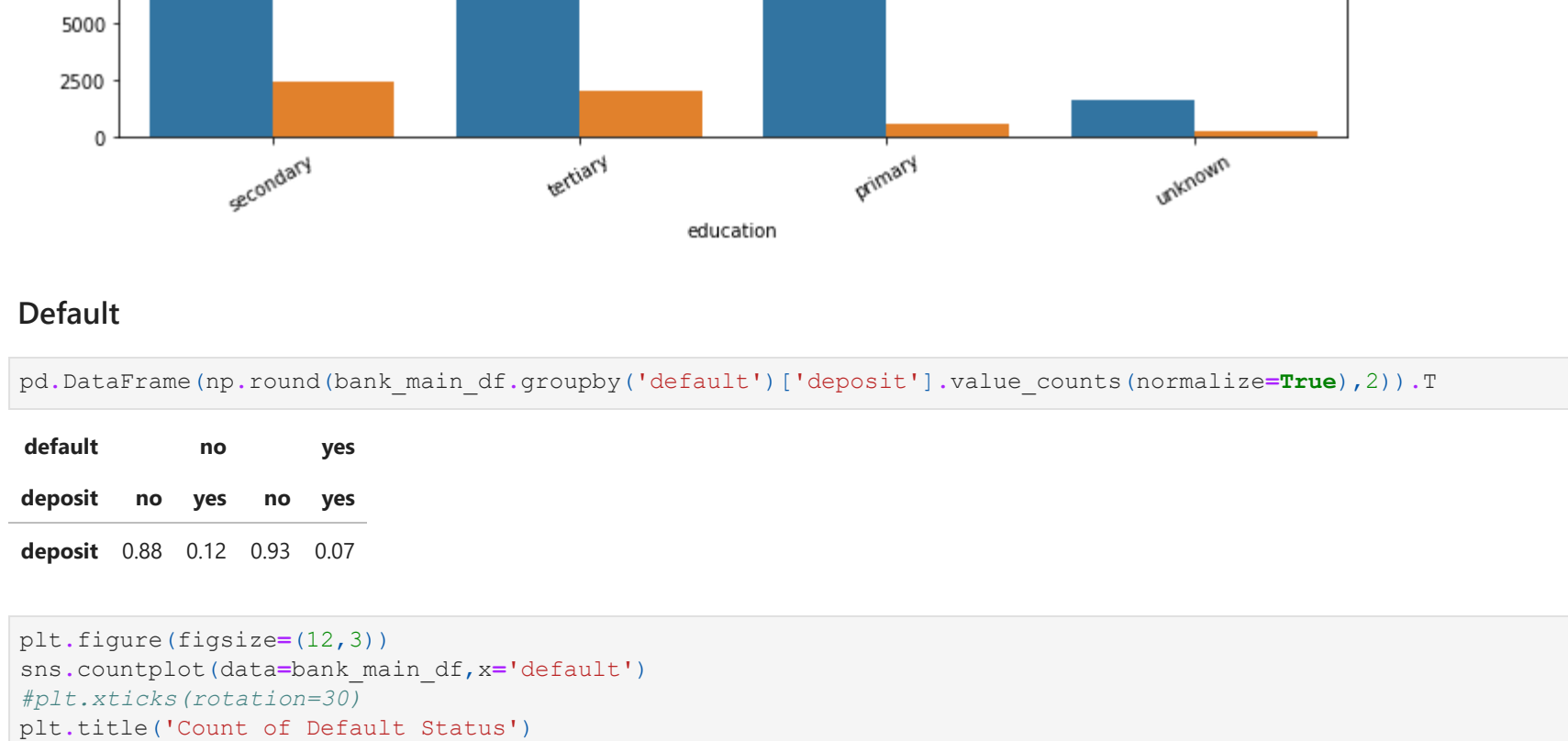
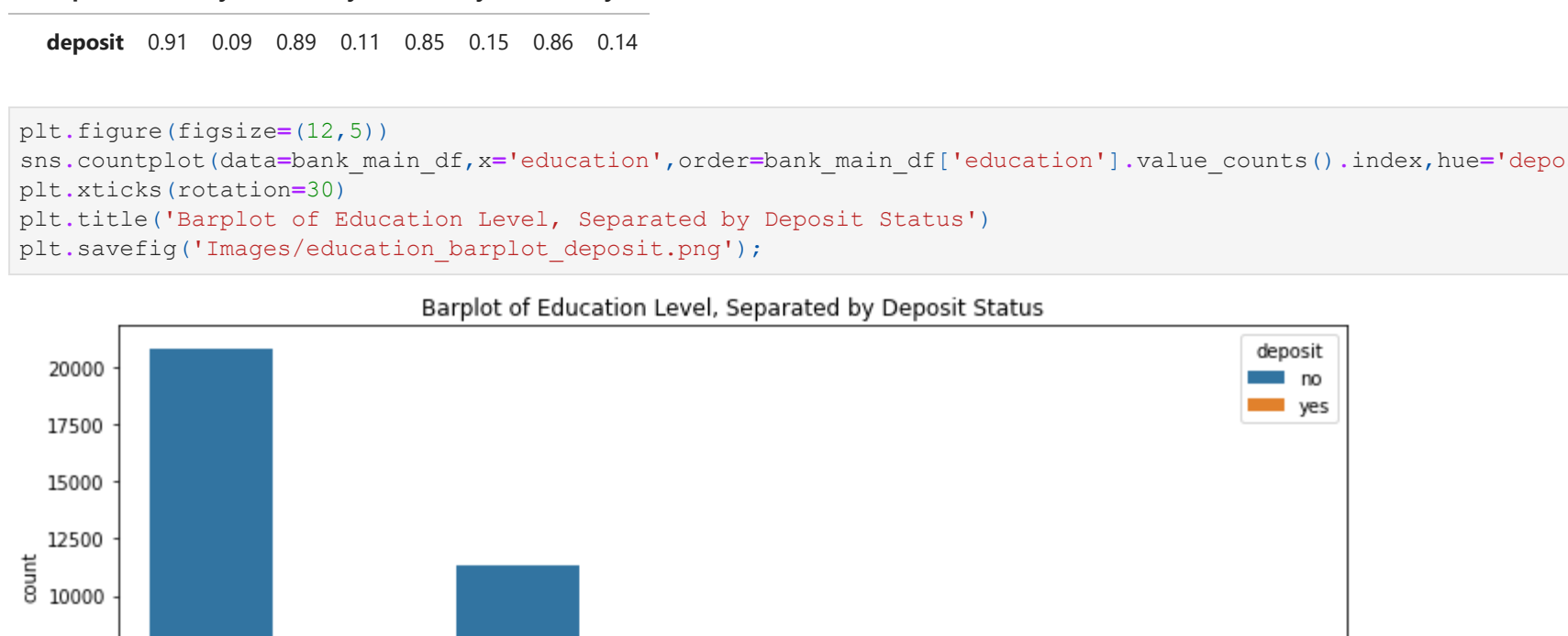
deposit
```



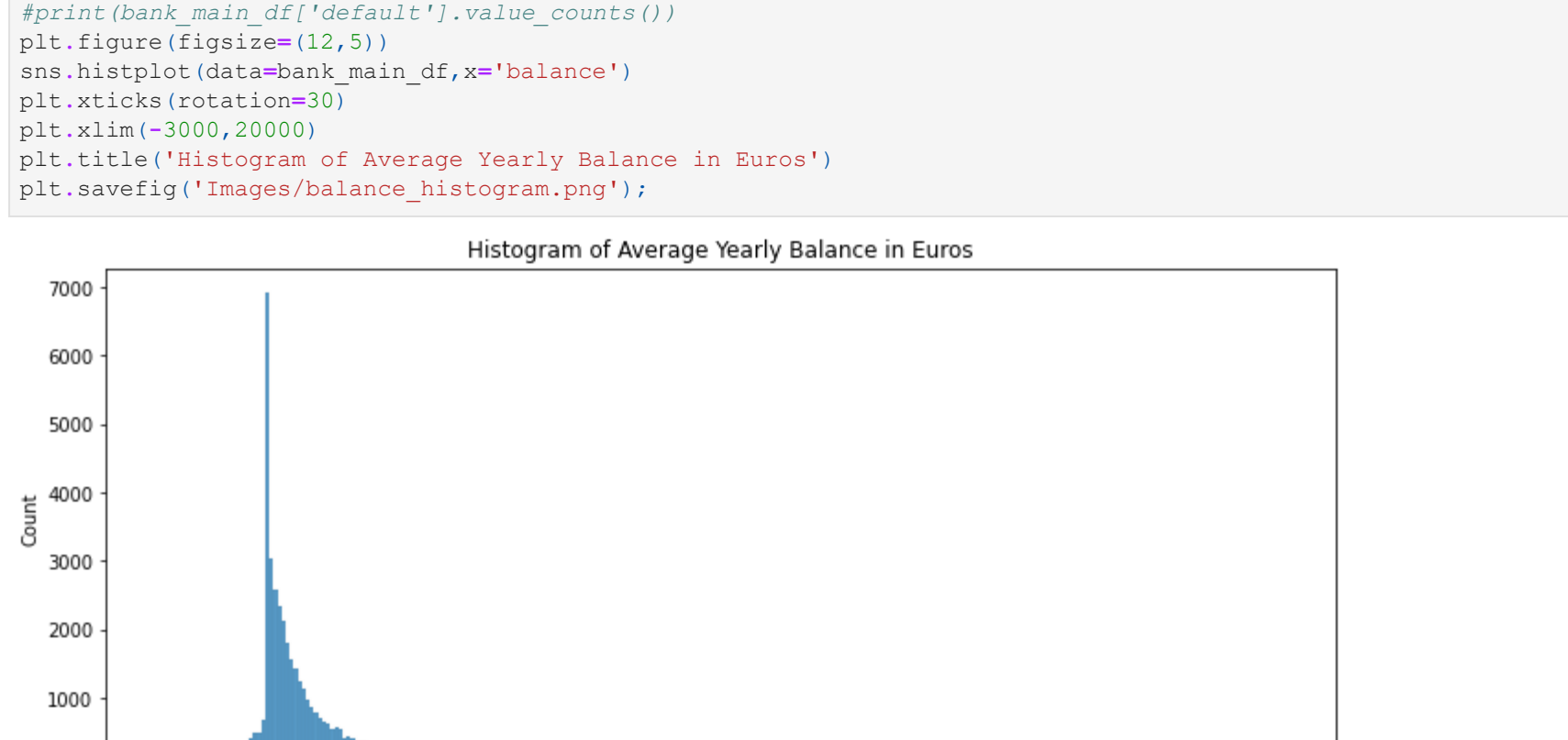
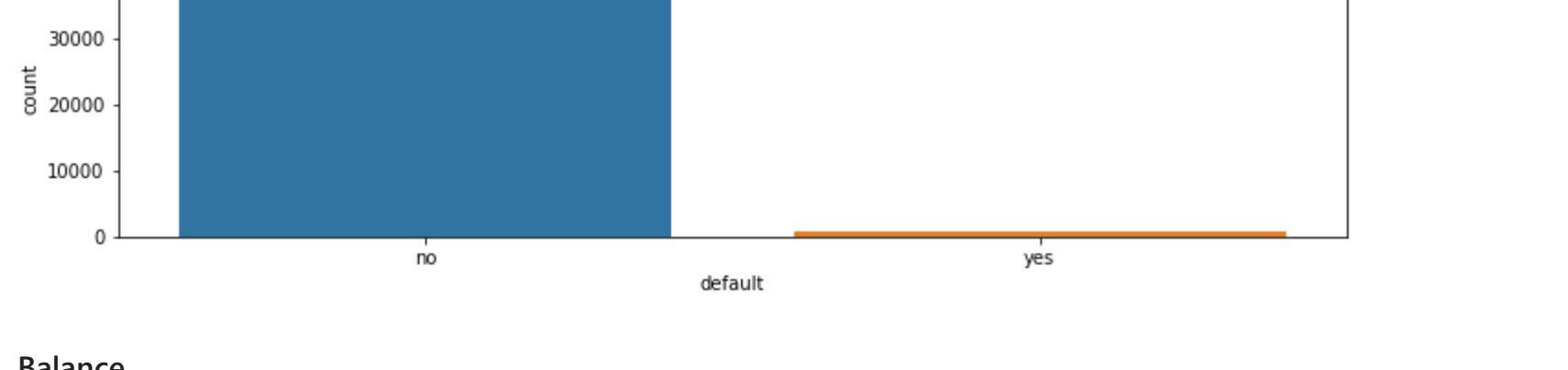
Age



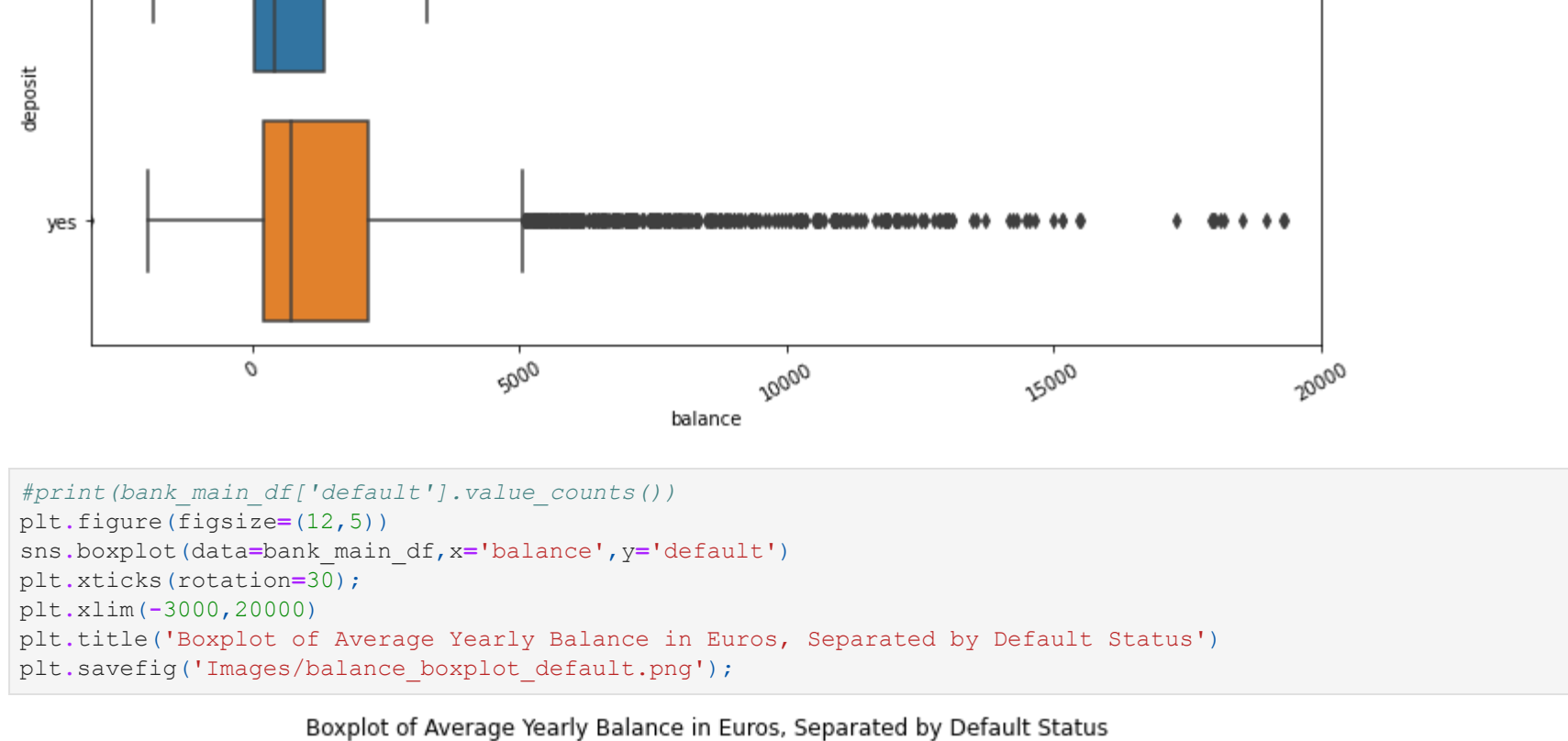
Job



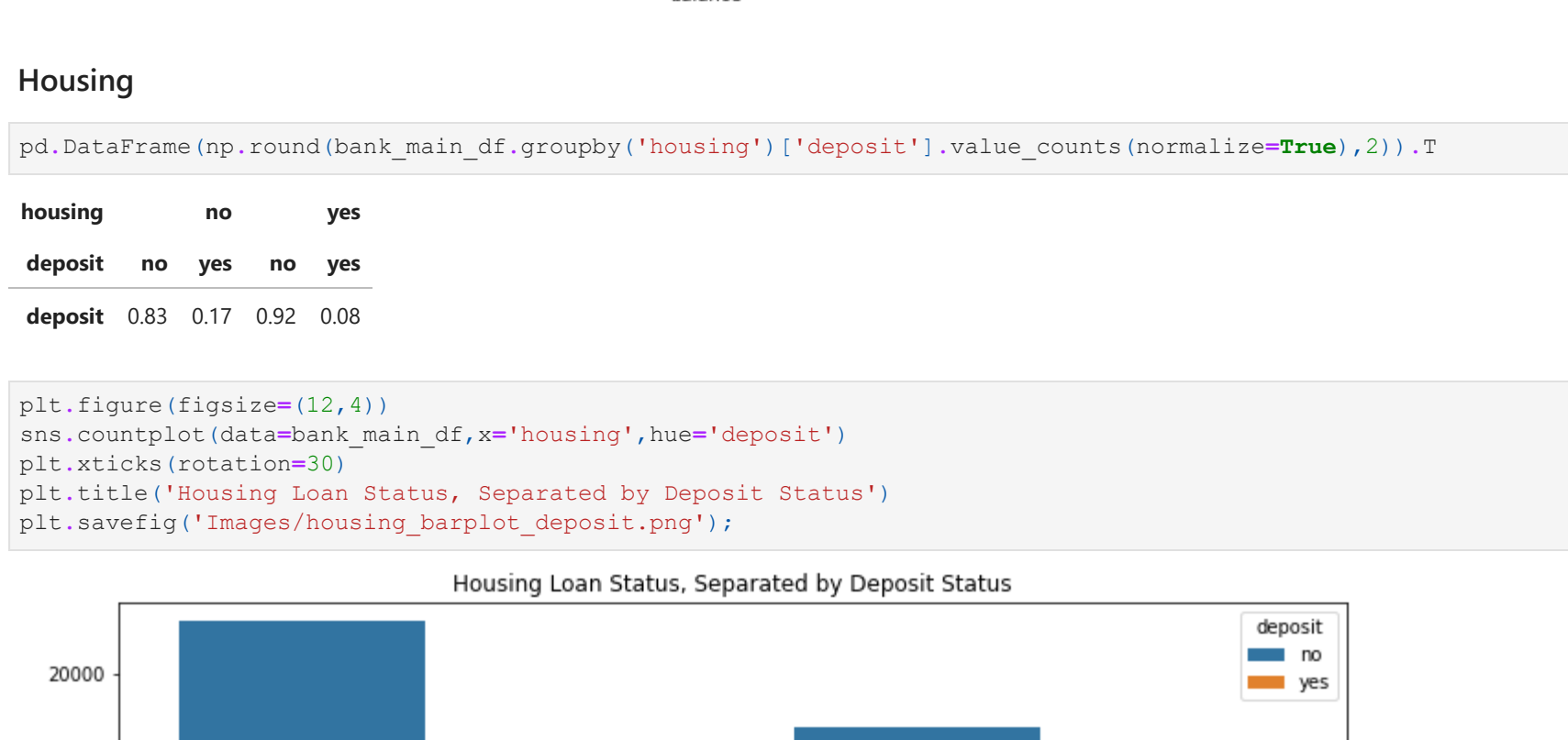
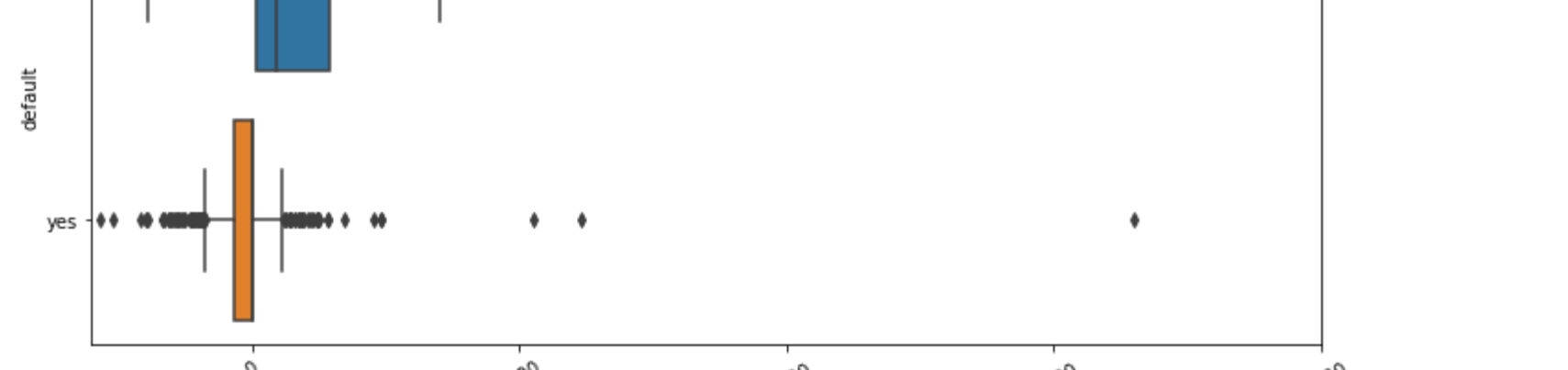
Marital



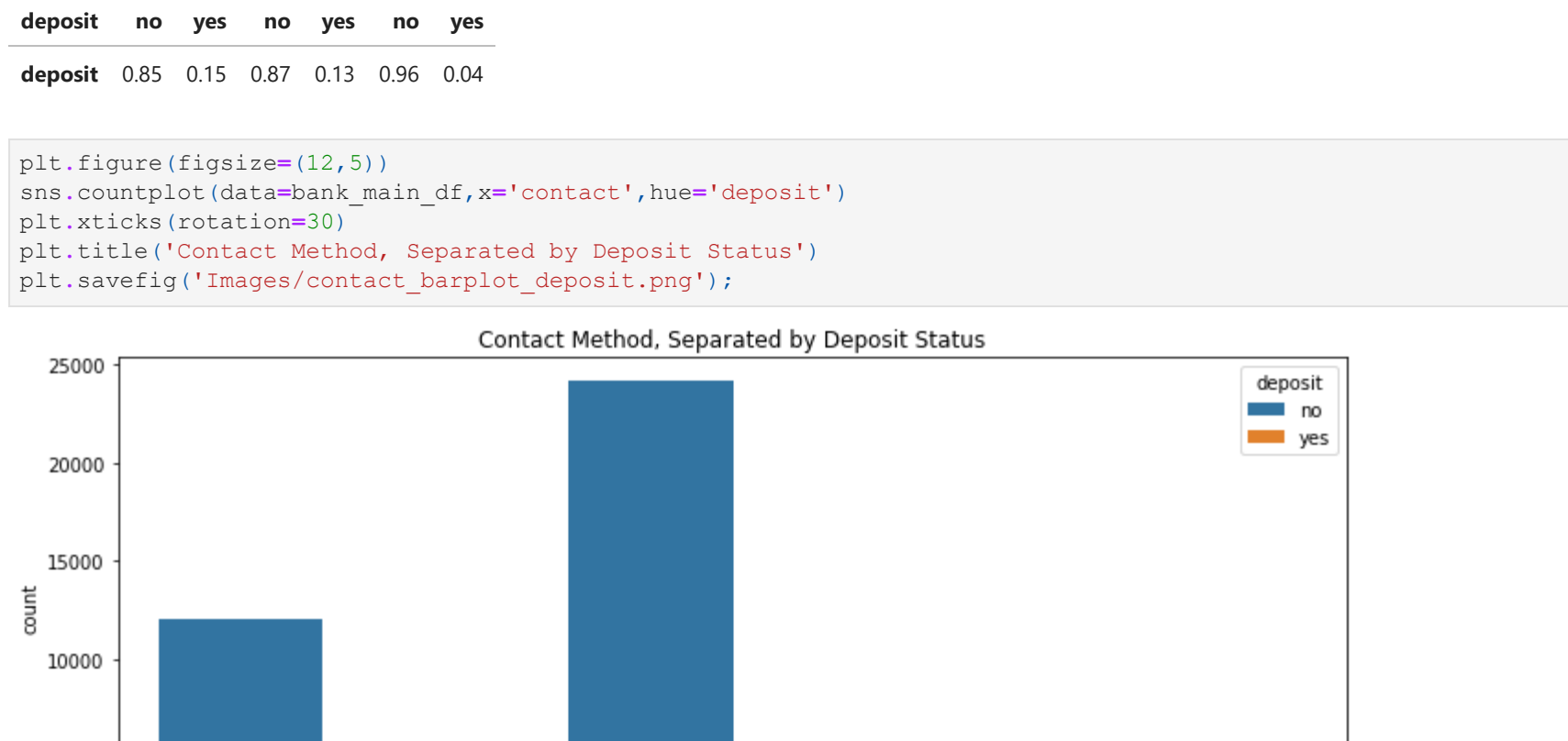
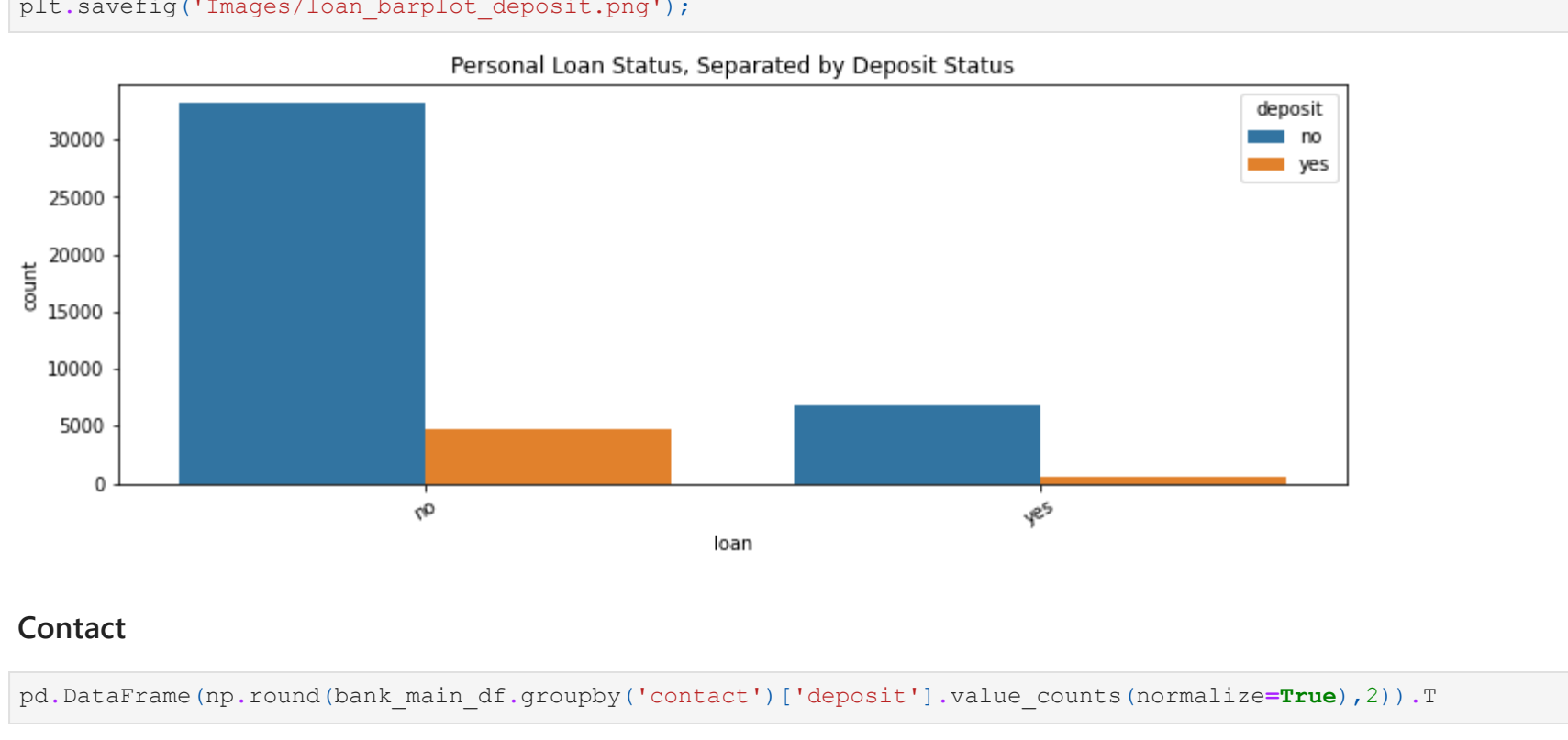
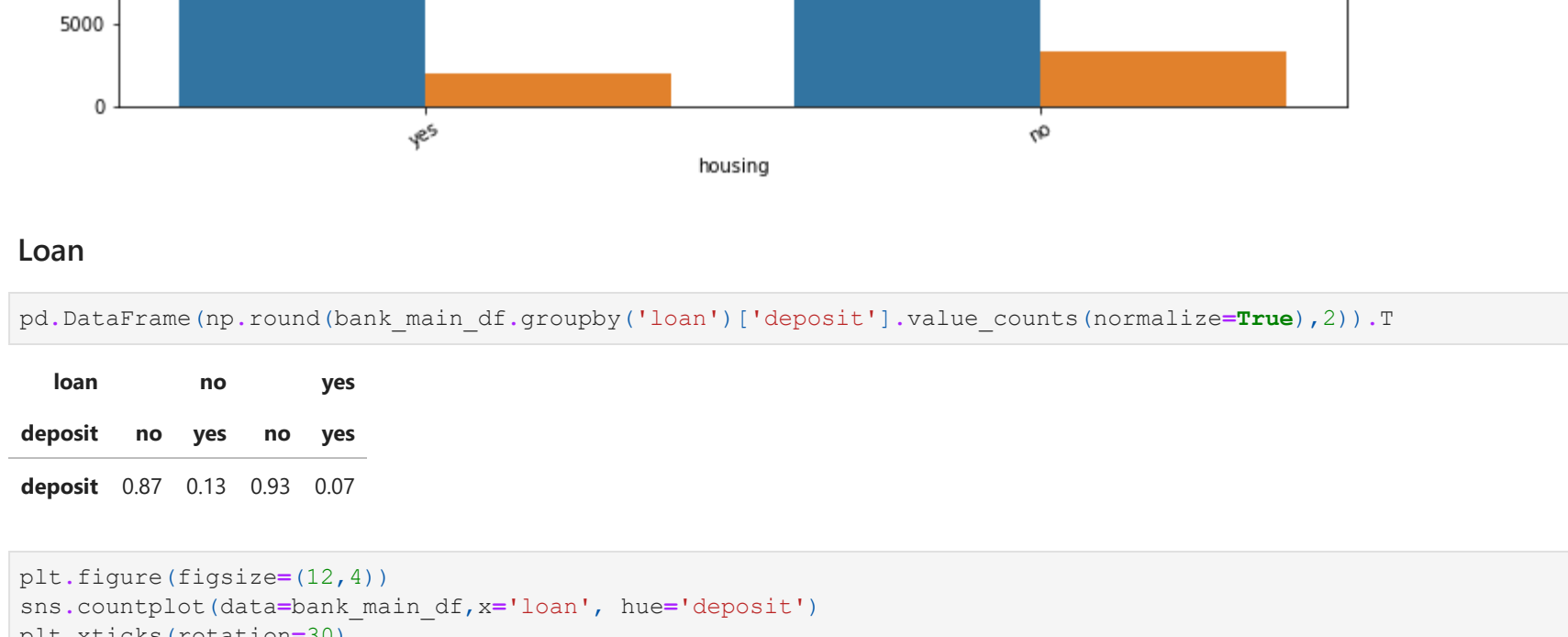
Education



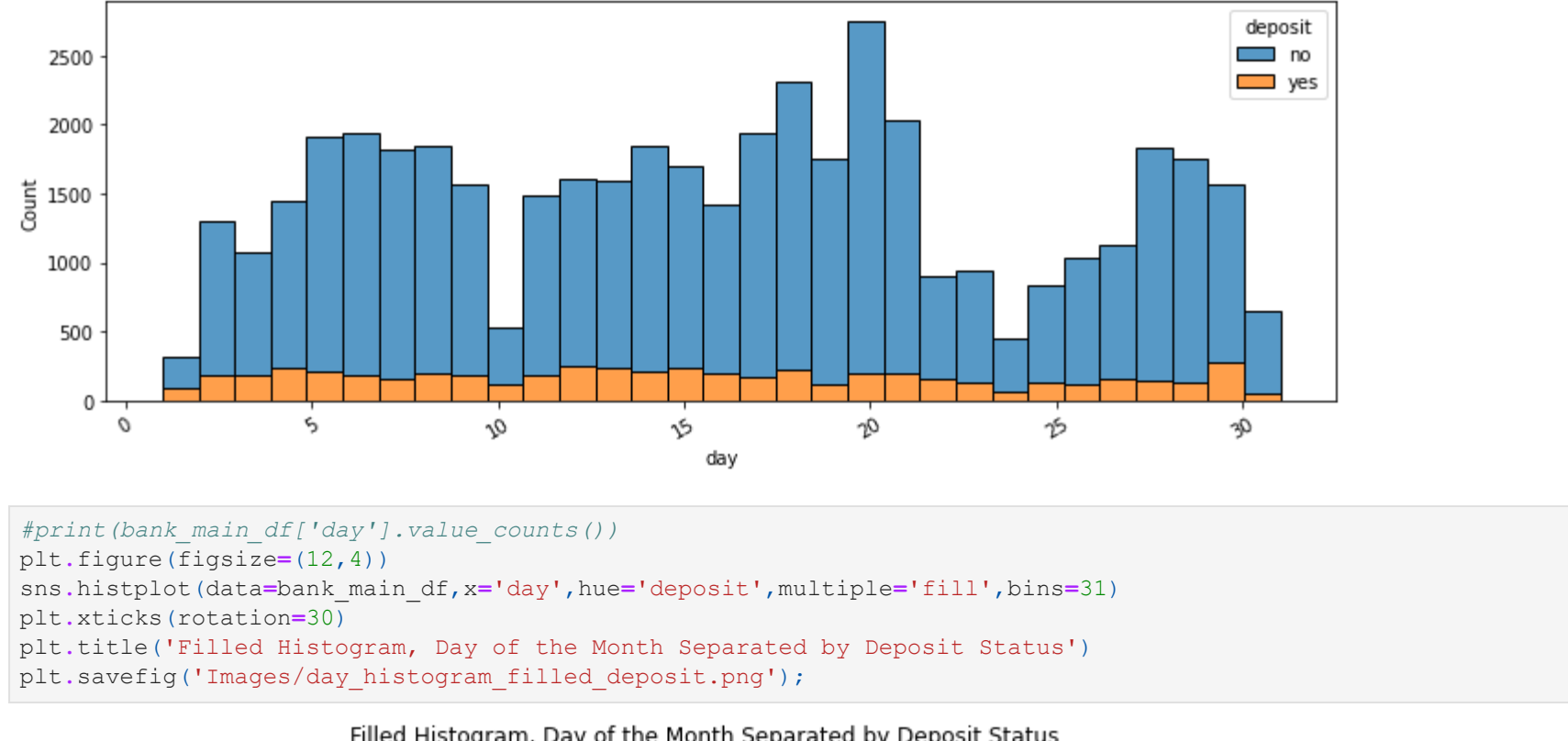
Default



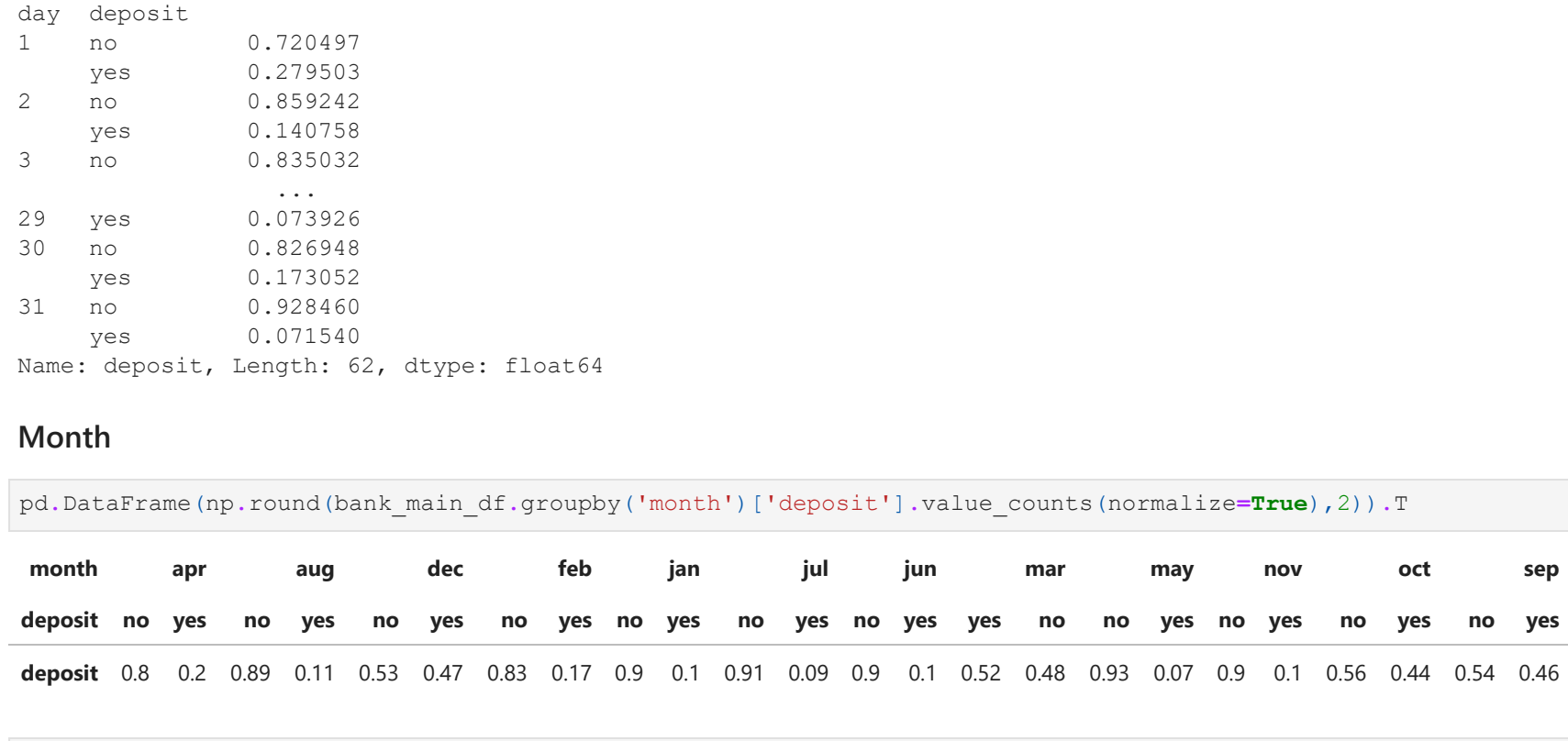
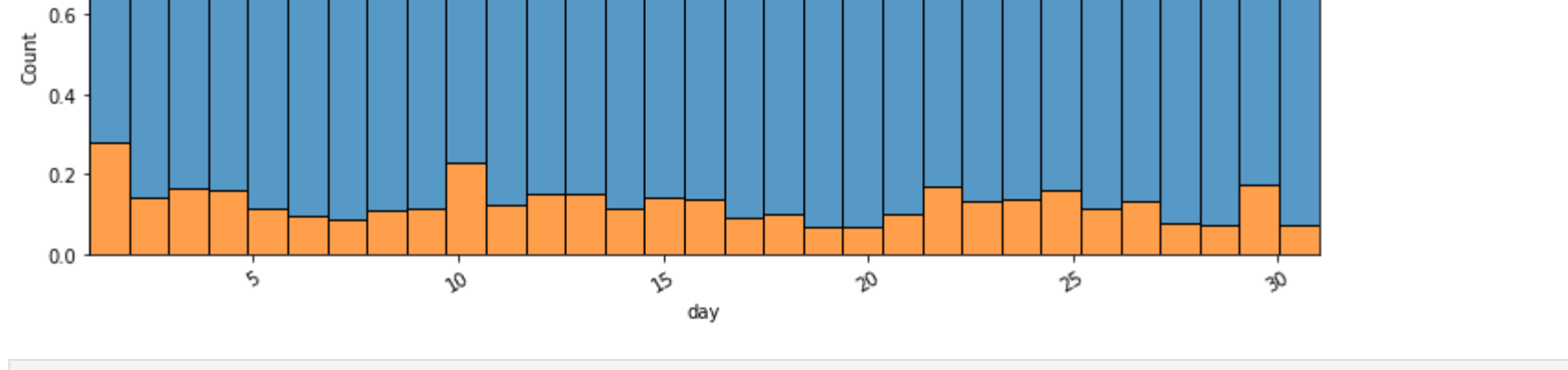
Balance



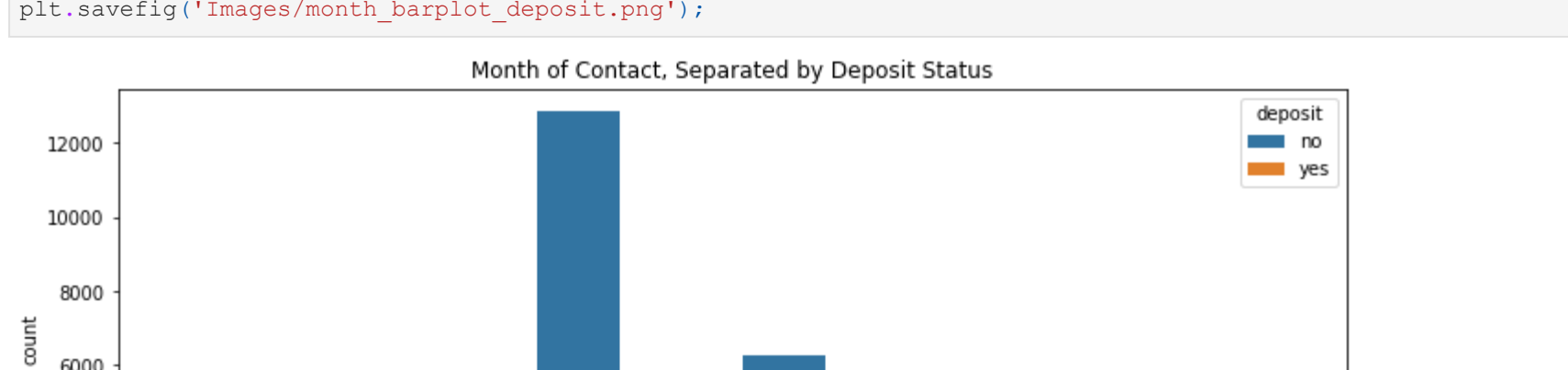
Housing



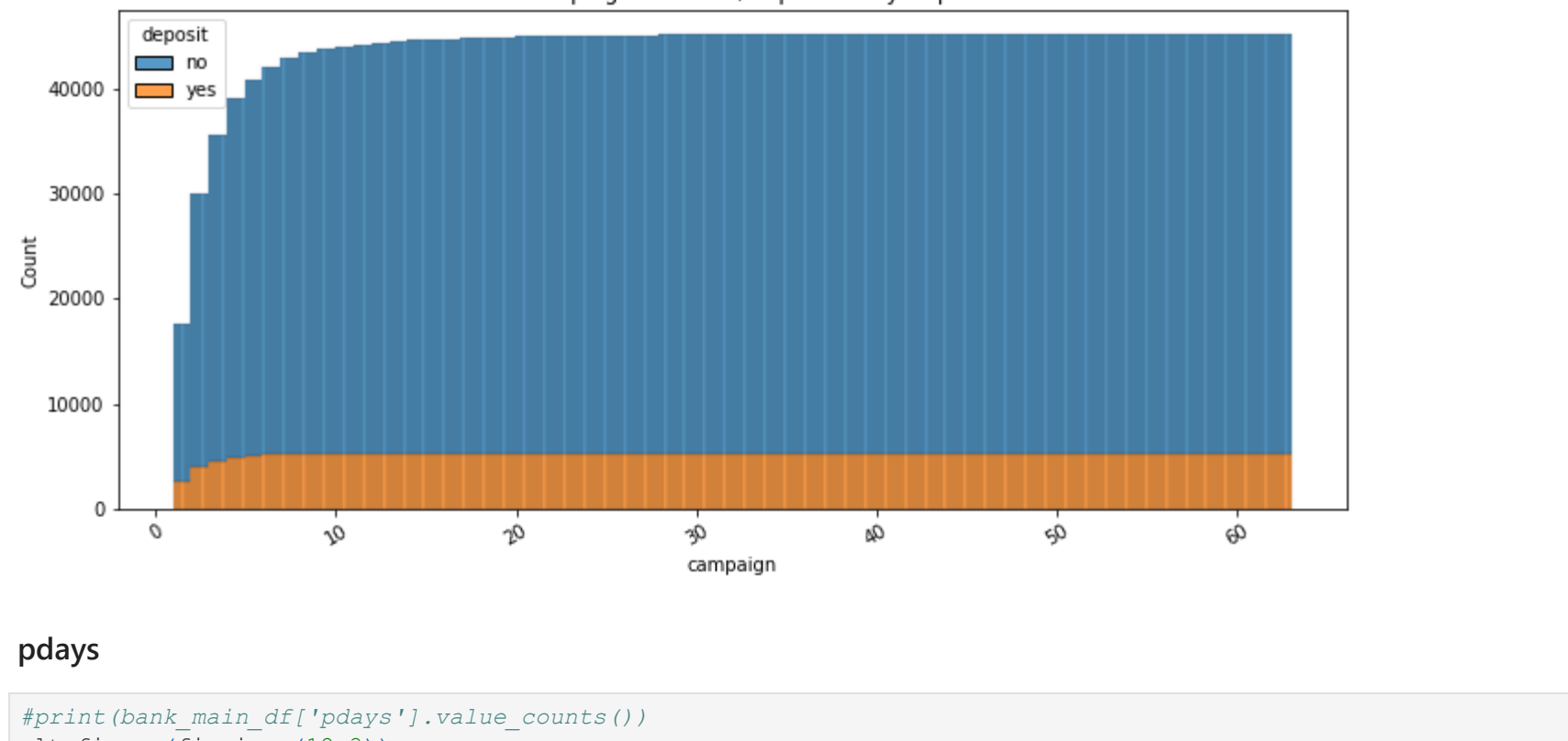
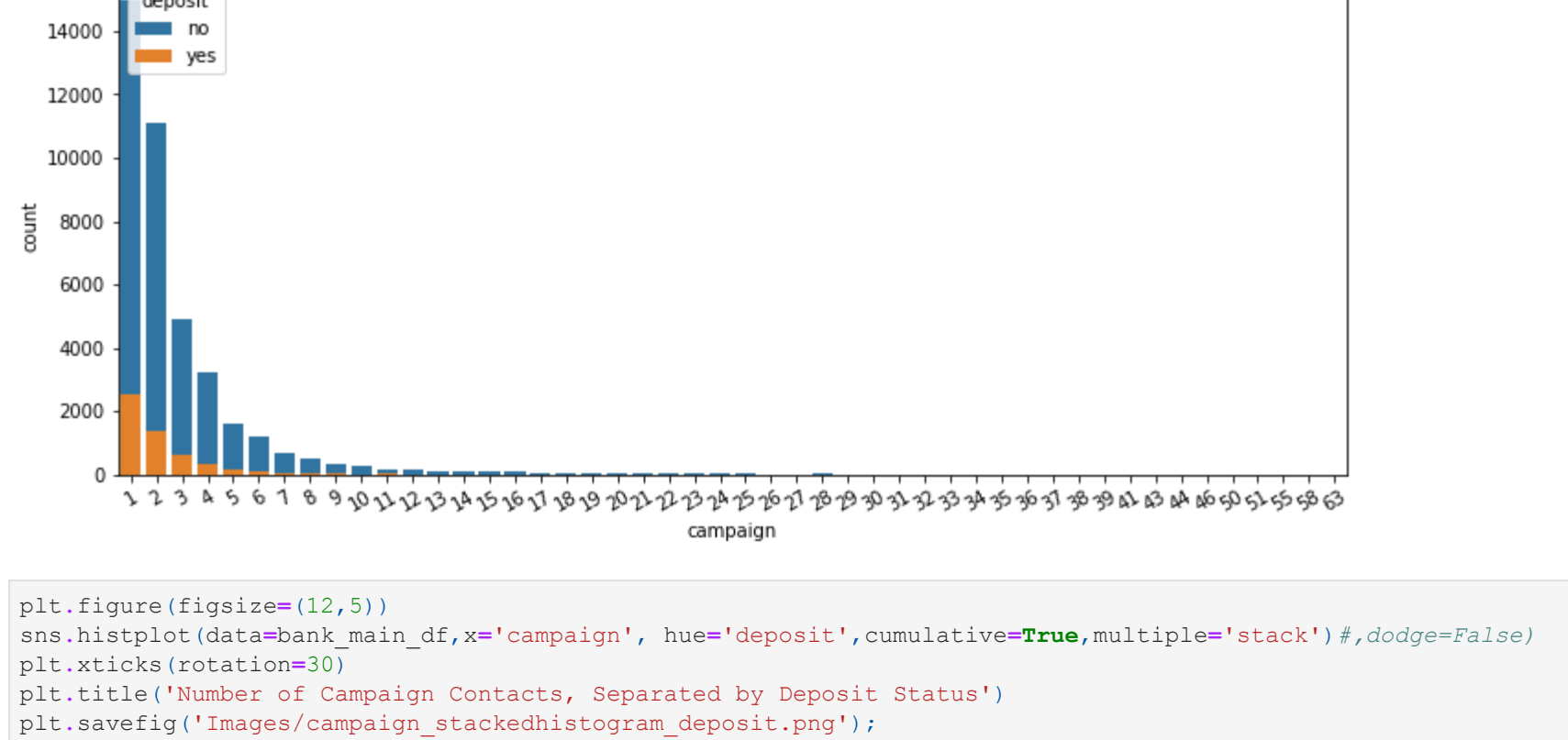
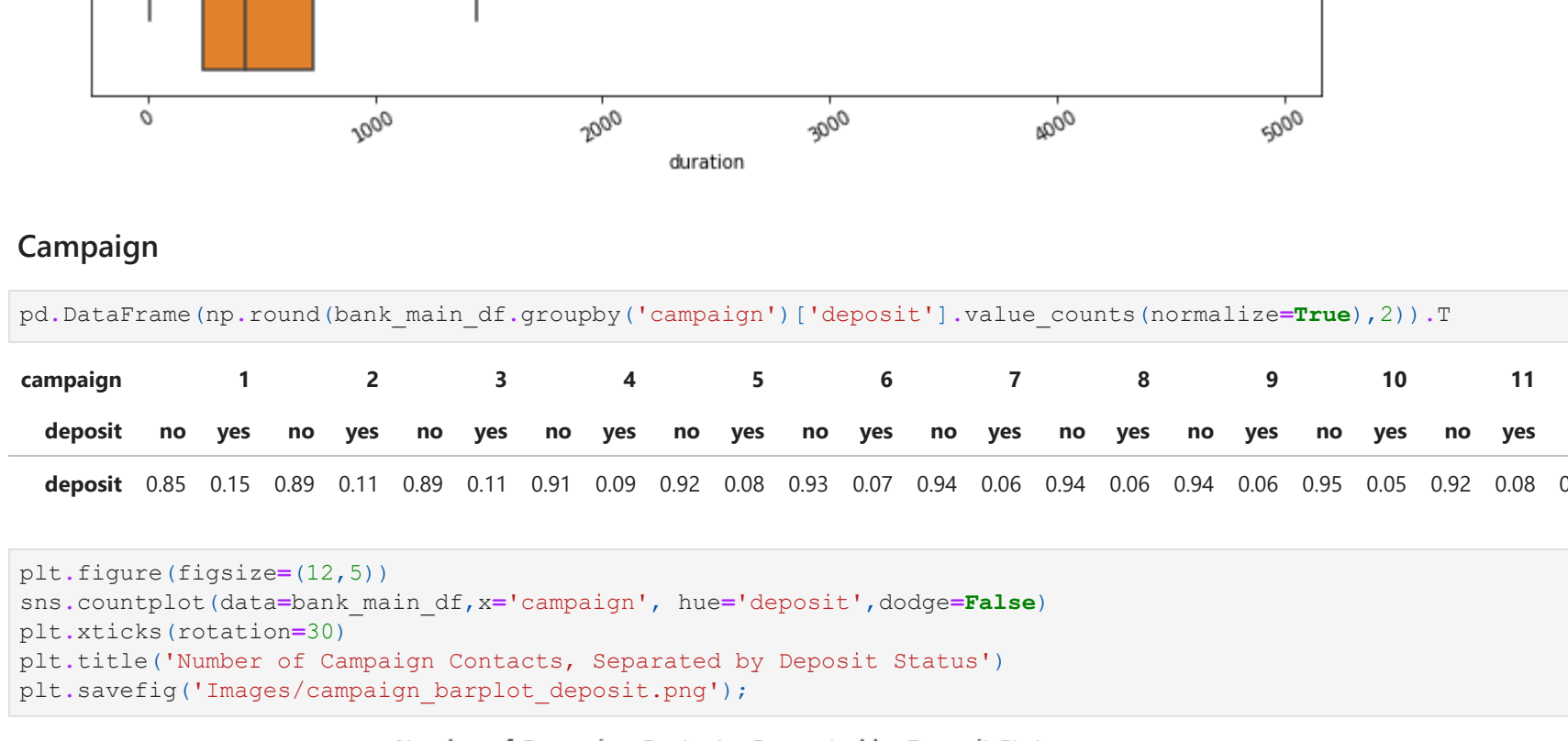
Loan



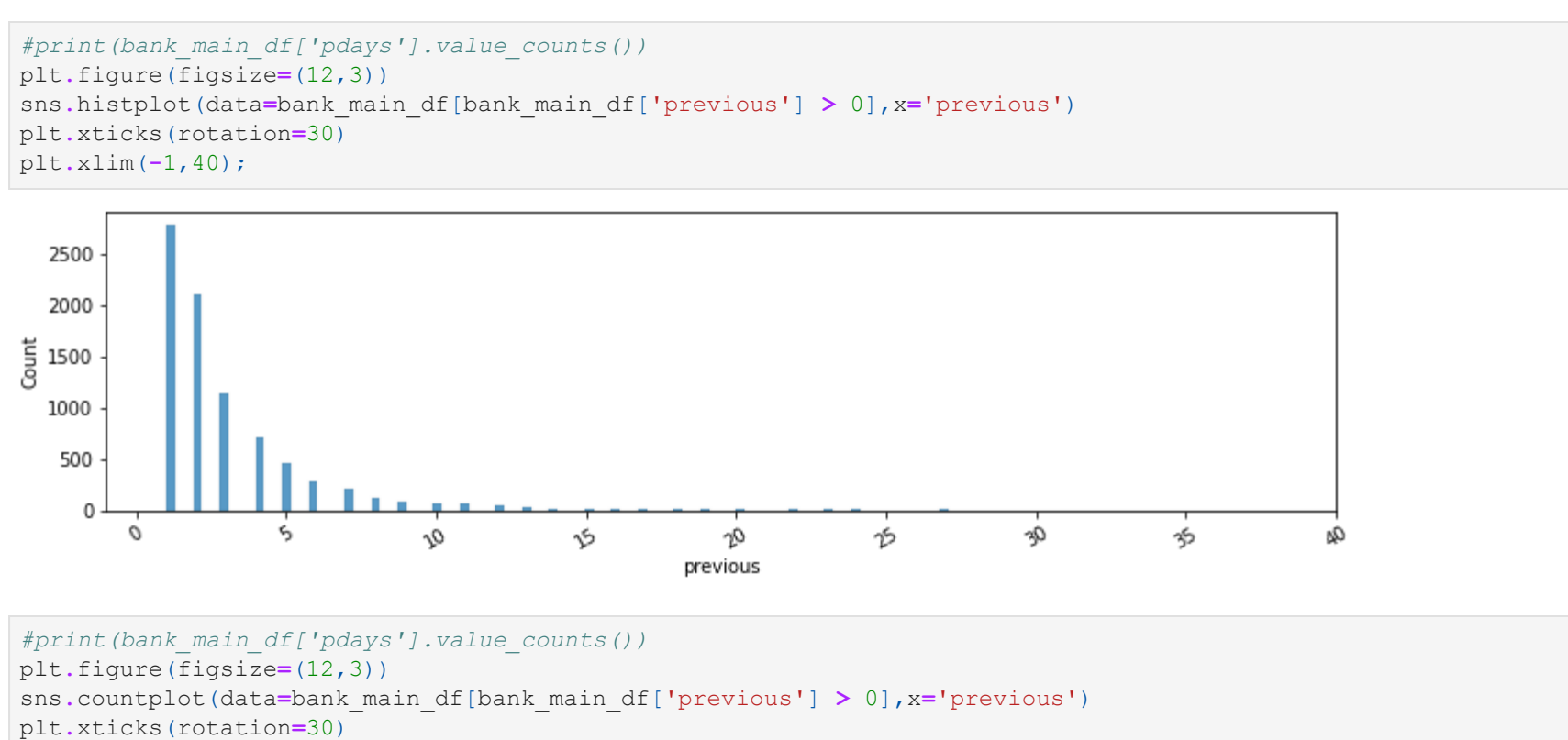
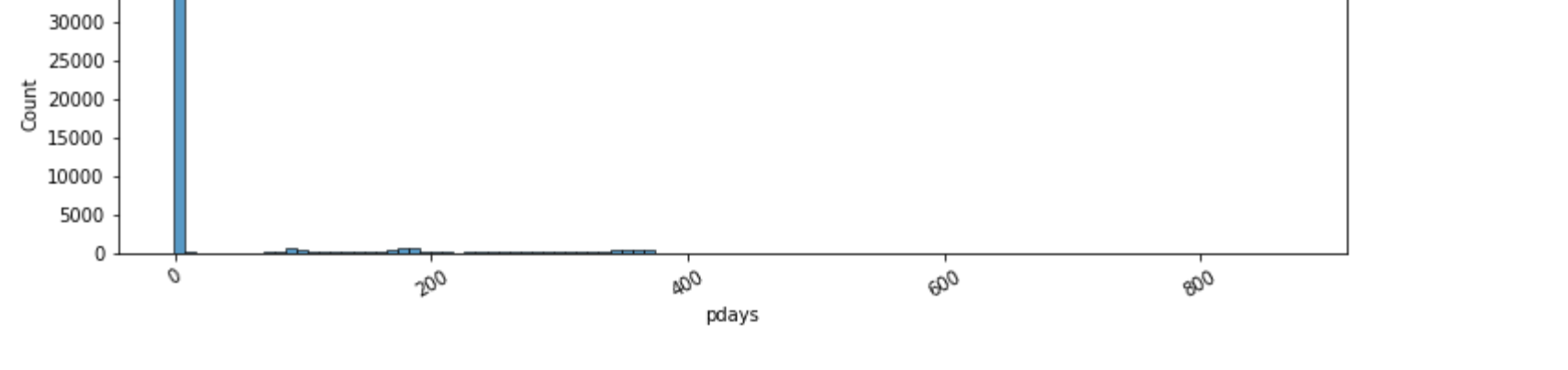
Contact



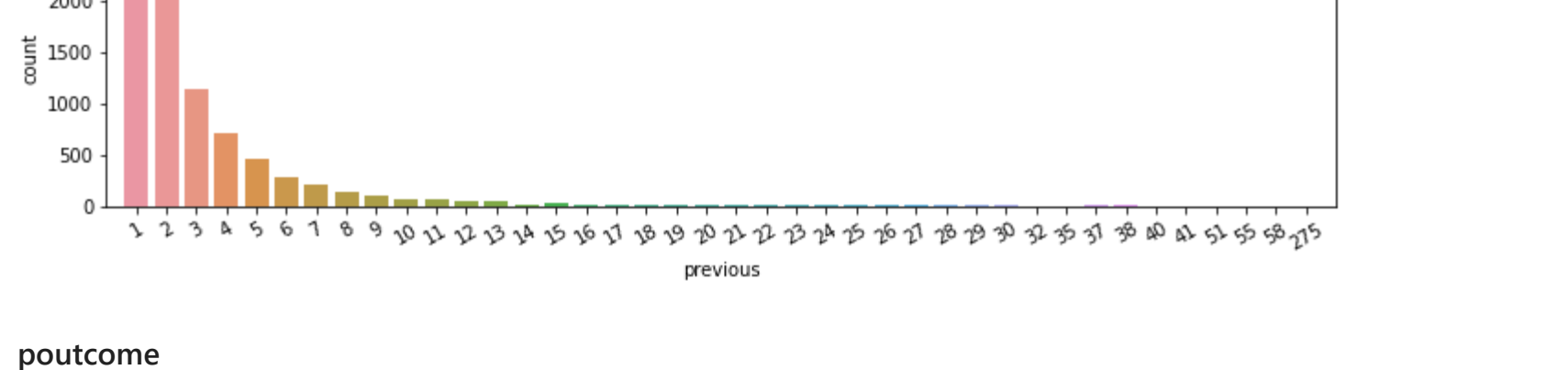
Day



Duration



Campaign

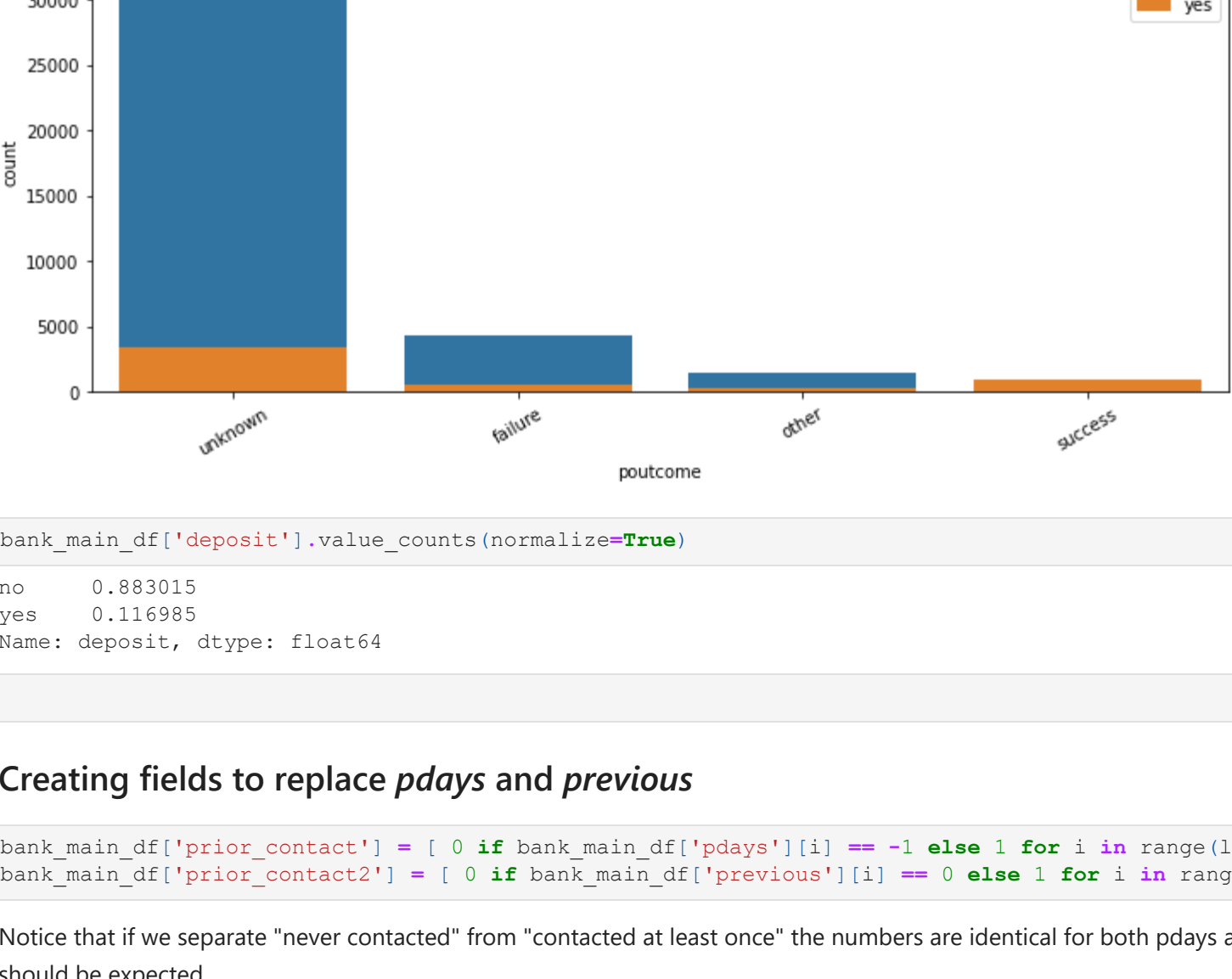


pdays

previous

poutcome


```
unknown      36959
failure       4901
other         1840
success       1511
Name: poutcome, dtype: int64
```



```
In [43]: bank_main_df['deposit'].value_counts(normalize=True)
```

```
Out[43]: no      0.883015
         yes     0.116985
         Name: deposit, dtype: float64
```

```
In [ ]:
```

Creating fields to replace *pdays* and *previous*

```
In [44]: bank_main_df['prior_contact'] = [ 0 if bank_main_df['pdays'][i] == -1 else 1 for i in range(len(bank_main_df))]
         bank_main_df['prior_contact2'] = [ 0 if bank_main_df['previous'][i] == 0 else 1 for i in range(len(bank_main_df))]
```

Notice that if we separate "never contacted" from "contacted at least once" the numbers are identical for both *pdays* and *previous*, as should be expected.

```
In [45]: plt.figure(figsize=(12,4))
sns.countplot(data=bank_main_df,hue='deposit',x='prior_contact',dodge=False)
plt.title("Was the Client Ever Contacted for Another Campaign, Yes (1) or No (0)")
plt.savefig('Images/prior_contact_yesno.png')
```



```
In [46]: bank_main_df.groupby('deposit')['prior_contact'].value_counts()#normalize=True)
```

```
Out[46]: deposit prior_contact    count
         0              33570
         1              6352
         yes            3384
         1              1905
         Name: prior_contact, dtype: int64
```

```
In [47]: sns.countplot(data=bank_main_df,x='deposit',hue='prior_contact2')
```

```
Out[47]: <AxesSubplot: xlabel='deposit', ylabel='count'>
```



```
In [48]: bank_main_df.groupby('deposit')['prior_contact2'].value_counts()#normalize=True)
```

```
Out[48]: deposit prior_contact2    count
         0              33570
         1              6352
         yes            3384
         1              1905
         Name: prior_contact2, dtype: int64
```

```
In [ ]:
```