

```
# Algebra
import numpy as np

# DataFrame
import pandas as pd

# Visualization
import matplotlib.pyplot as plt
import seaborn as sns

# Algorithms
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import OneHotEncoder
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import plot_confusion_matrix

from sklearn.model_selection import cross_validate
from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import StratifiedKFold
from sklearn.ensemble import RandomForestClassifier

import imblearn
import scipy
import math
```

Tasks to Perform :

Read in the file and get basic information about the data, including numerical summaries

```
In [2]: bdata = pd.read_csv(r'F:\ML\bank_marketing_data.csv')

In [3]: # NO of Rows and Columns in data :
bdata.shape

Out[3]: (45211, 19)

In [4]: # Checking the Presence of Null Values :
bdata.isnull().sum()

Out[4]:
age                0
job                0
salary            0
marital           0
education         0
targeted         0
default          0
balance          0
loan             0
contact          0
day              0
month            0
duration         0
campaign         0
pdays          0
previous         0
poutcome        0
response        0
dtype: int64
```

- Observation : No null values present in data

```
In [5]: # Checking DataTypes in Data :
bdata.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 45211 entries, 0 to 45210
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   age                    45211 non-null    int64
1   job                    45211 non-null    object
2   salary                 45211 non-null    int64
3   marital                45211 non-null    object
4   education              45211 non-null    object
5   targeted               45211 non-null    object
6   default                45211 non-null    object
7   balance                45211 non-null    int64
8   housing                45211 non-null    object
9   loan                   45211 non-null    object
10  contact                45211 non-null    object
11  day                     45211 non-null    int64
12  month                  45211 non-null    object
13  duration               45211 non-null    int64
14  campaign               45211 non-null    int64
15  pdays                  45211 non-null    int64
16  previous               45211 non-null    int64
17  poutcome               45211 non-null    object
18  response               45211 non-null    object
dtypes: int64(8), object(11)
memory usage: 6.6+ MB
```

- Observation : Two types of datatypes: int and object

```
In [6]: # Checking description of numerical class :
bdata.describe()

Out[6]:
```

	age	salary	balance	day	duration	campaign	pdays
count	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000	45211.000000
mean	40.936210	57006.171065	1362.272058	15.806419	258.163080	2.763841	40.197828
std	10.618762	32085.718415	3044.765829	8.322476	257.527812	3.098021	100.128746
min	18.000000	0.000000	-8019.000000	1.000000	0.000000	1.000000	-1.000000
25%	33.000000	20000.000000	72.000000	8.000000	103.000000	1.000000	-1.000000
50%	39.000000	60000.000000	448.000000	16.000000	319.000000	2.000000	-1.000000
75%	48.000000	70000.000000	1428.000000	21.000000	180.000000	3.000000	-1.000000
max	95.000000	120000.000000	102127.000000	31.000000	4918.000000	63.000000	871.000000

- Observation :

1. Age of customer is Almost Normally Distributed

1. Salary of Customer is almost Normally Distributed

1. Account balance is skewed

```
In [7]: # Checking the Skewness of Numerical columns :
bdata.skew()

Out[7]:
age          0.684818
salary       0.137829
balance      8.360308
day          0.093079
duration     3.144318
campaign     4.488450
pdays       2.615715
previous     41.846454
dtype: float64
```

- Observation : balance is Skewed as has been described previously

- Describe pdays column make note of the mean , median and minimum values . Anything fishy in the values ?
- ANS : Pdays Number of days that passed by after the client was last contacted from a previous campaign

```
In [8]: bdata['pdays'].describe()

Out[8]:
count      45211.000000
mean       40.197828
std        100.128746
min        -1.000000
25%        -1.000000
50%        -1.000000
75%        -1.000000
max        871.000000
Name: pdays, dtype: float64
```

```
In [9]: bdata['pdays'].skew()

Out[9]: 2.6157154736563477

• Mean = 40.2 days(approx)

• Median = -1 and minimum = -1

• Yes there is fishy values in the data i.e. (-1), since no of days can't be NEGATIVE
```

Q. Describe the pdays column again , this time limiting yourself to relevant values of pdays.How different are mean and median values ?

```
In [10]: bdata_1 = bdata.loc[bdata['pdays'] > -1]

In [11]: bdata_1['pdays'].describe()

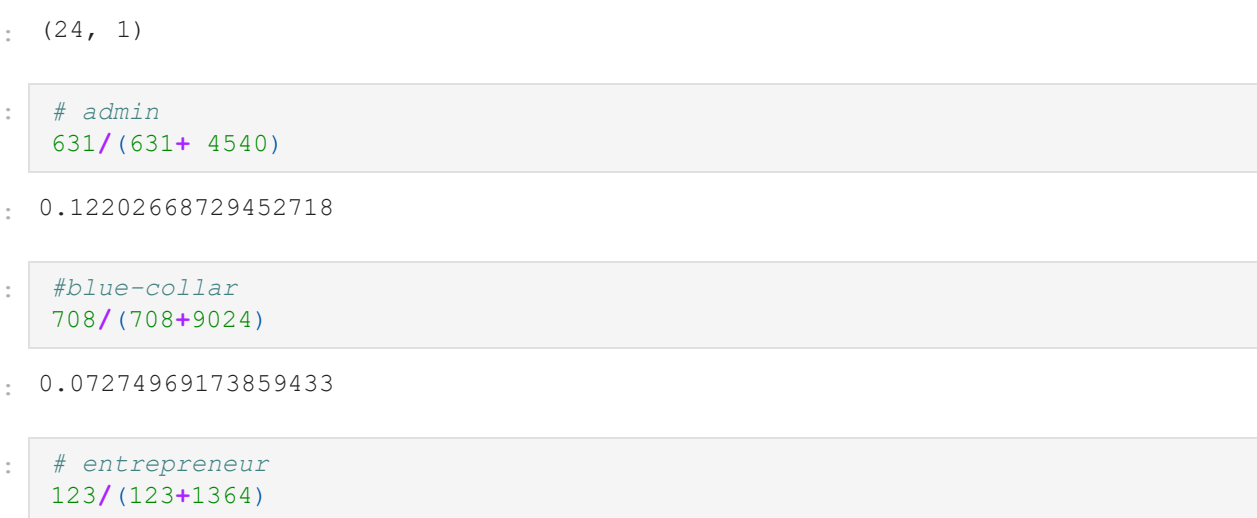
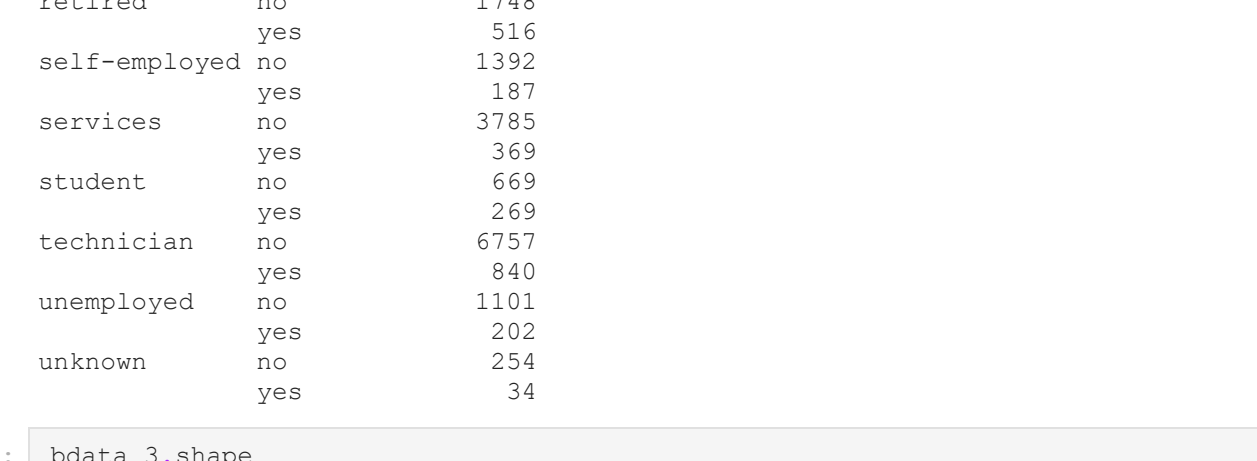
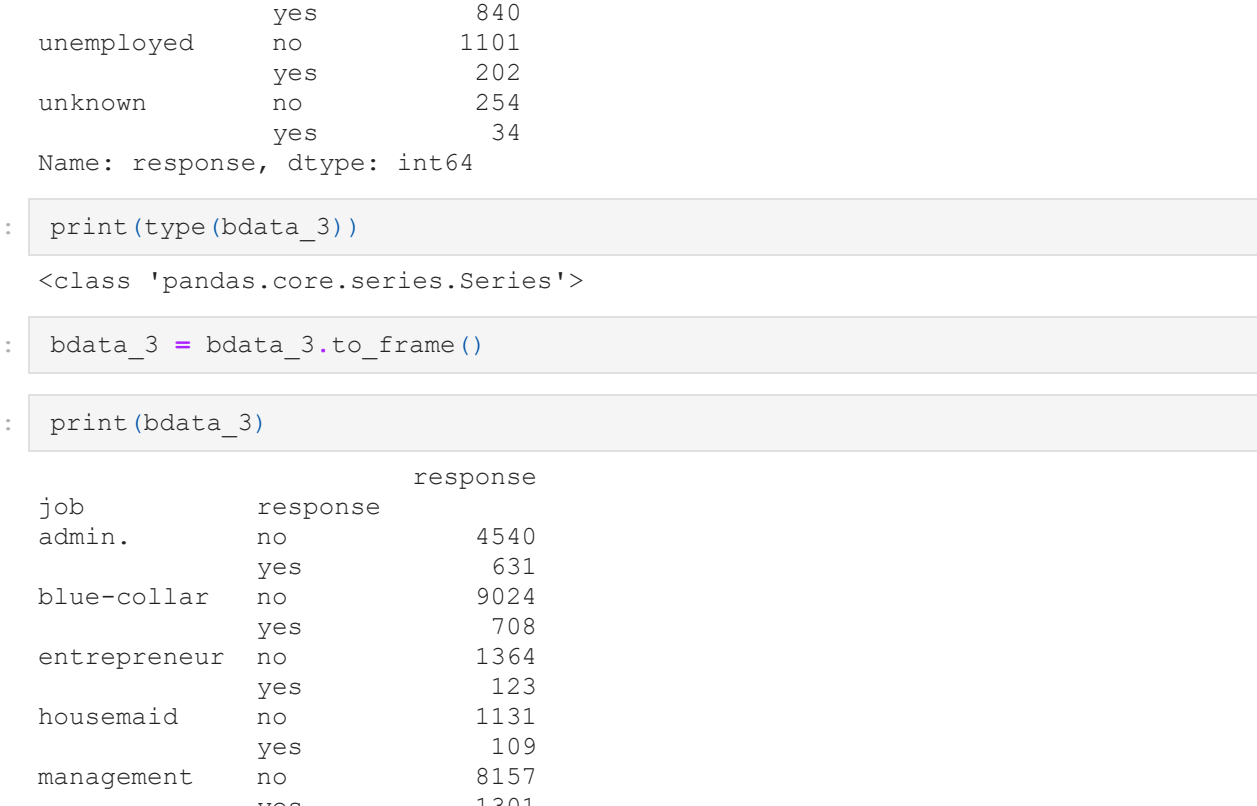
Out[11]:
count      8257.000000
mean       224.577692
std        115.344035
min         1.000000
25%        133.000000
50%        194.000000
75%        327.000000
max        871.000000
Name: pdays, dtype: float64
```

```
In [12]: bdata_1['pdays'].skew()

Out[12]: 0.6931397093928039

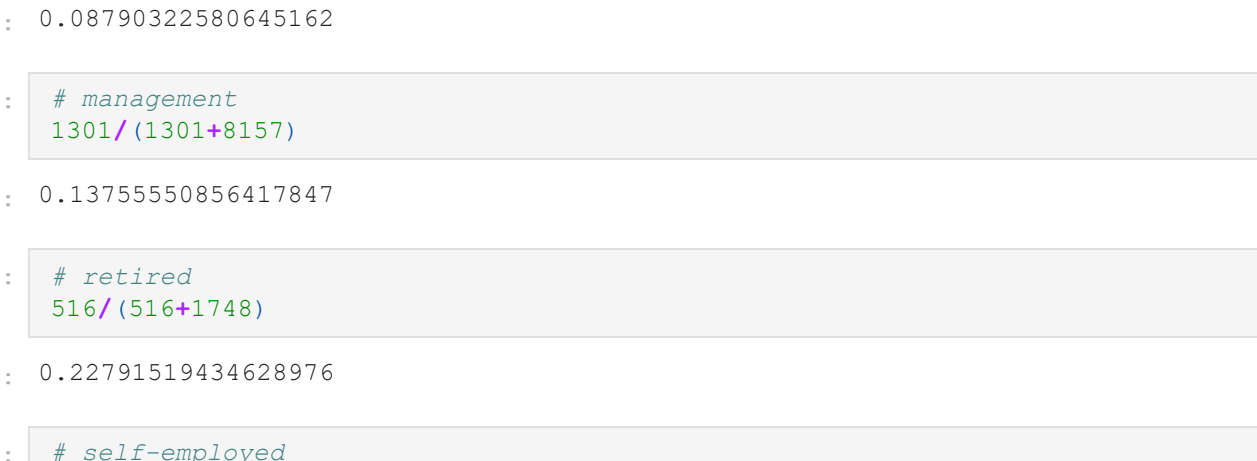
• Observation : previously mean was 40 days and median was -1 day but now it changed to 224 days and 194 days respectively
```

Q. Plot a horizontal bar graph with the median values of balance for each education level value. Which group has the highest median?



ANS: tertiary educational group has highest median

Q. Make a box plot for pdays. Do you see any outliers?



- Checking outliers after removing -1 values

ANS :yes, there is presence of outlier (days greater than 600 are outliers)

EDA Part :

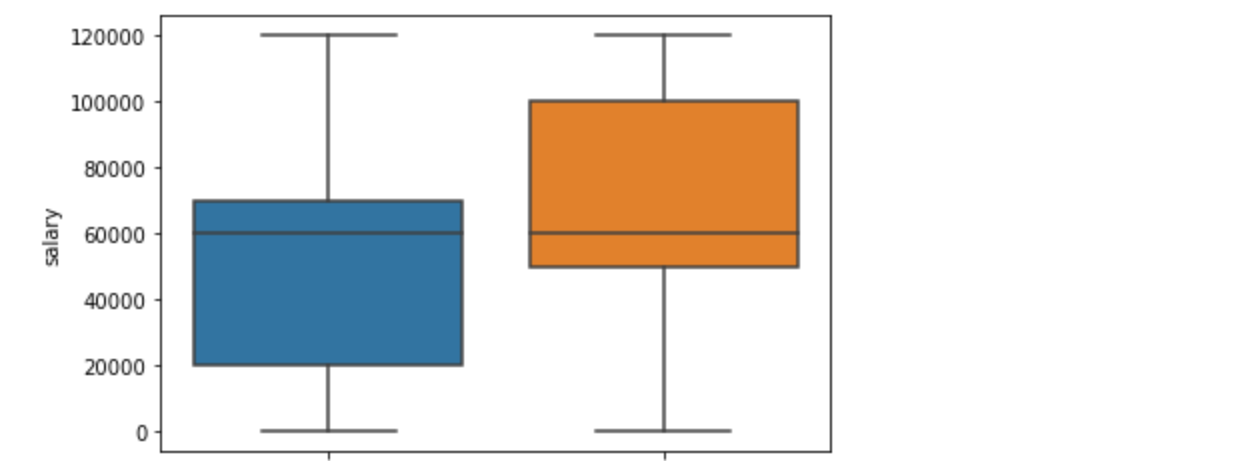
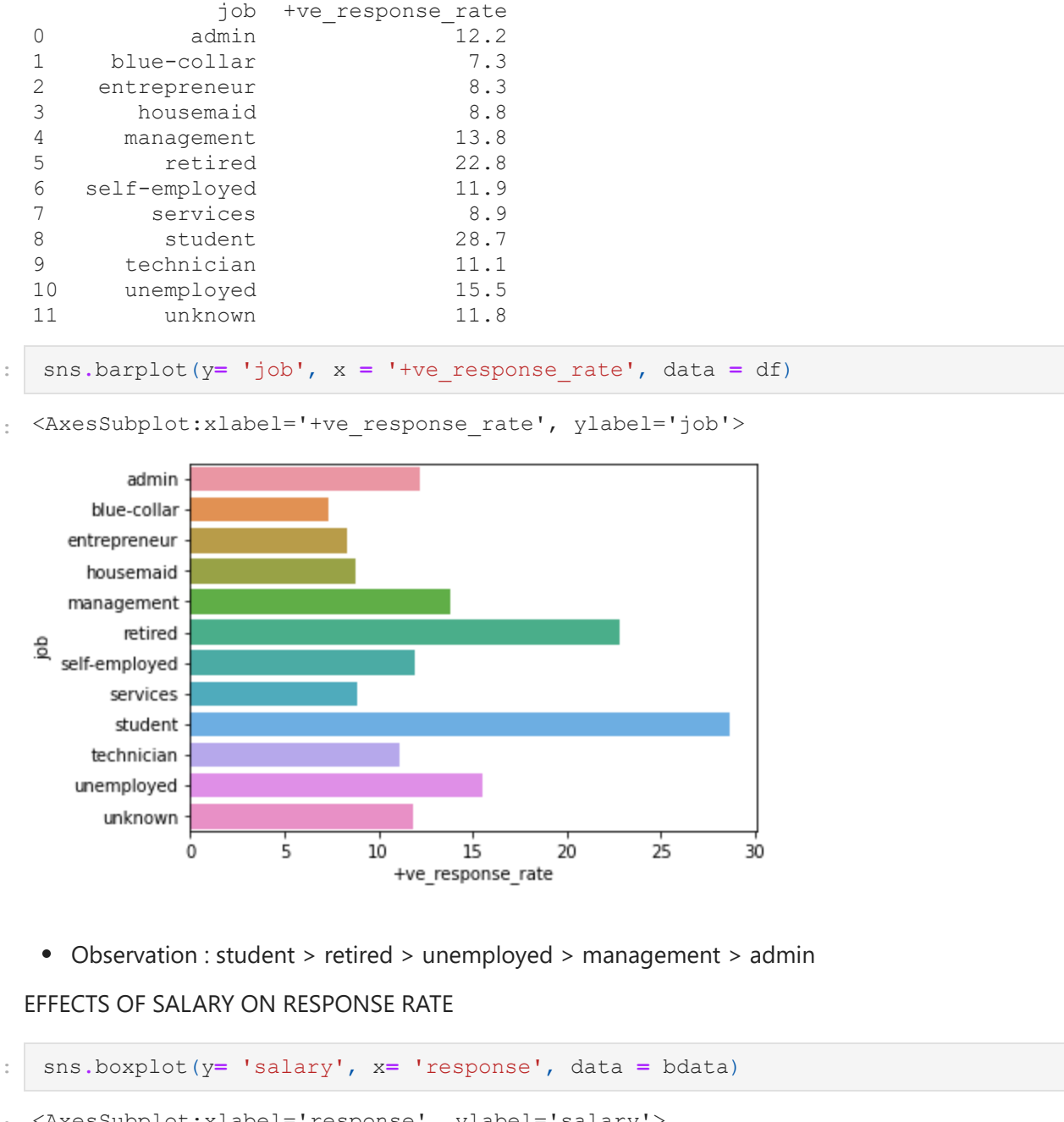
TARGET VARIABLE : response

- effects of AGE group on response



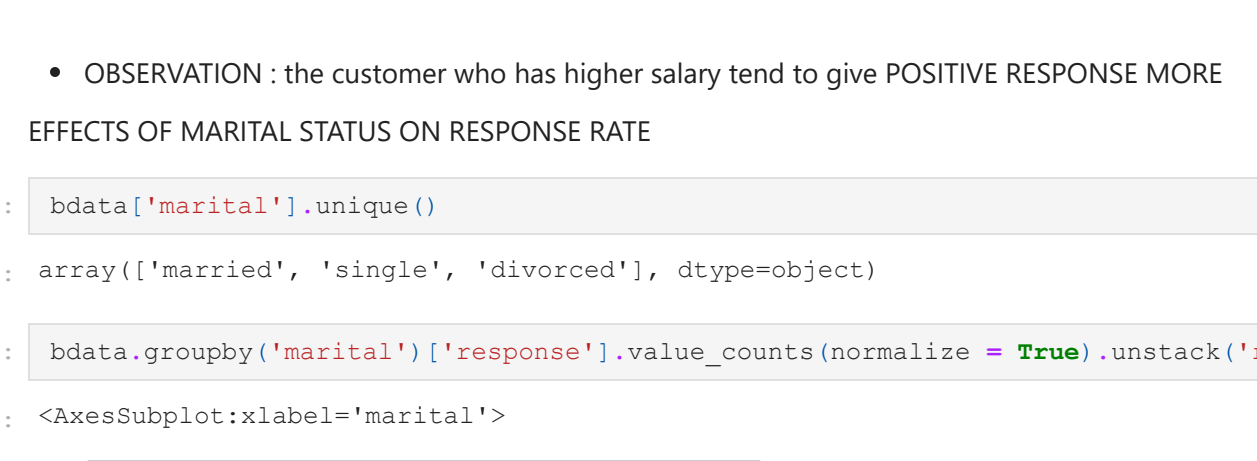
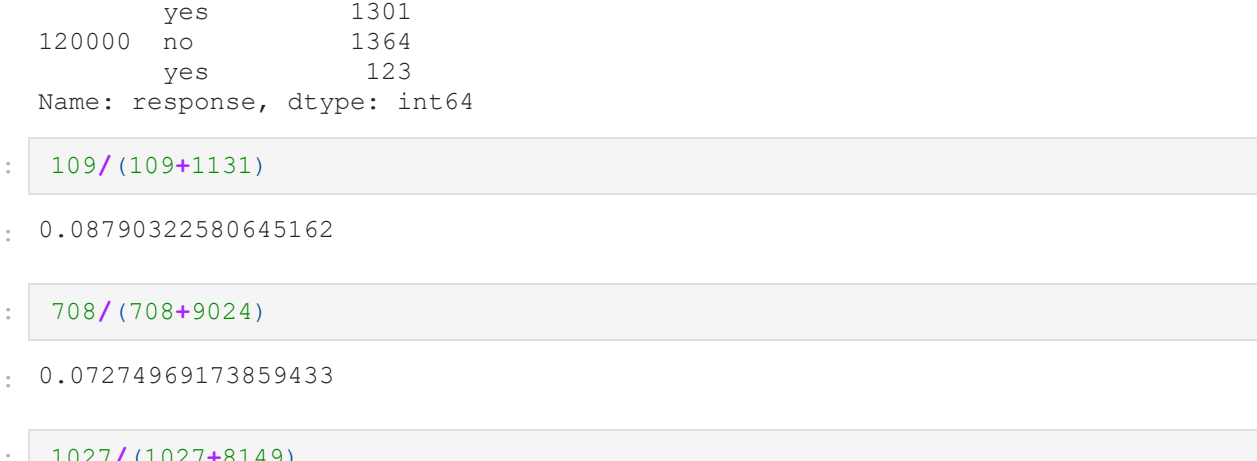
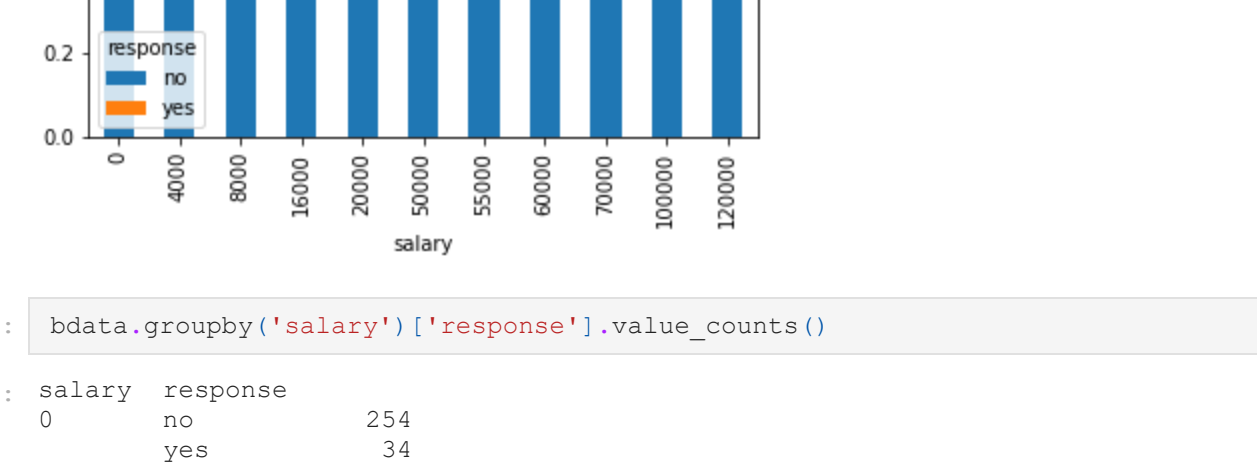
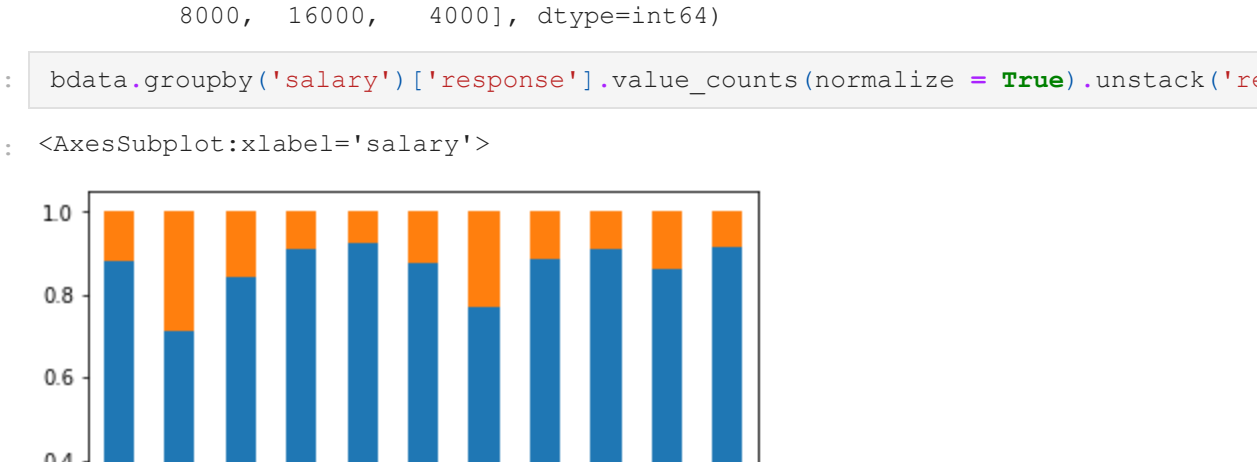
- OBSERVATION : THE age group doesn't have much effect on response of customer

EFFECTS OF JOB ON RESPONSE OF A CUSTOMER

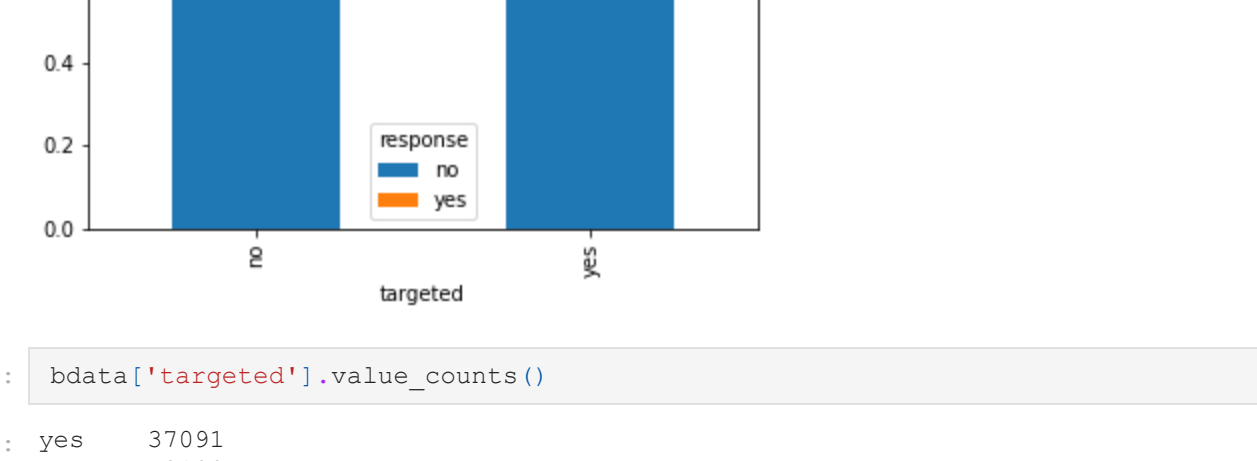
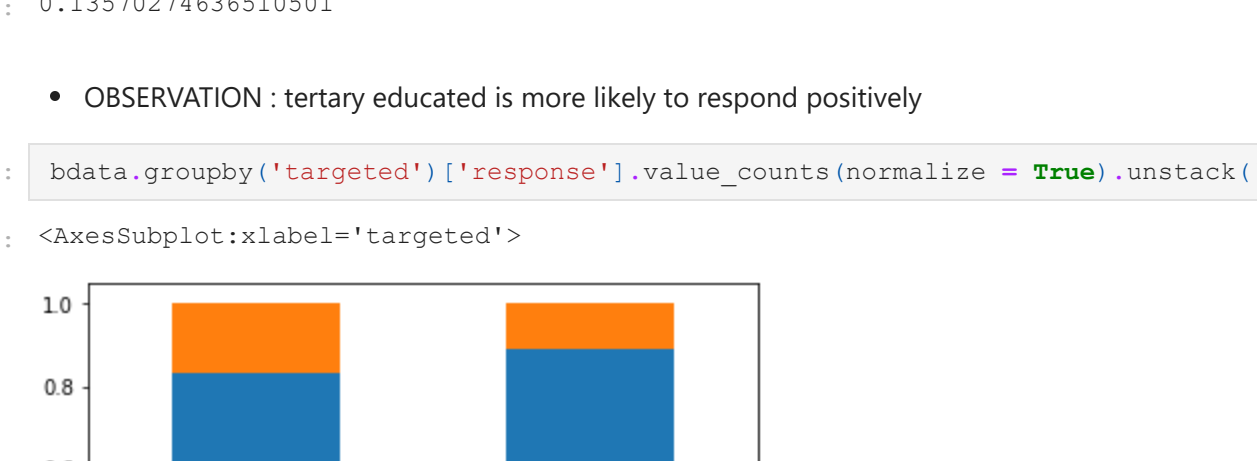
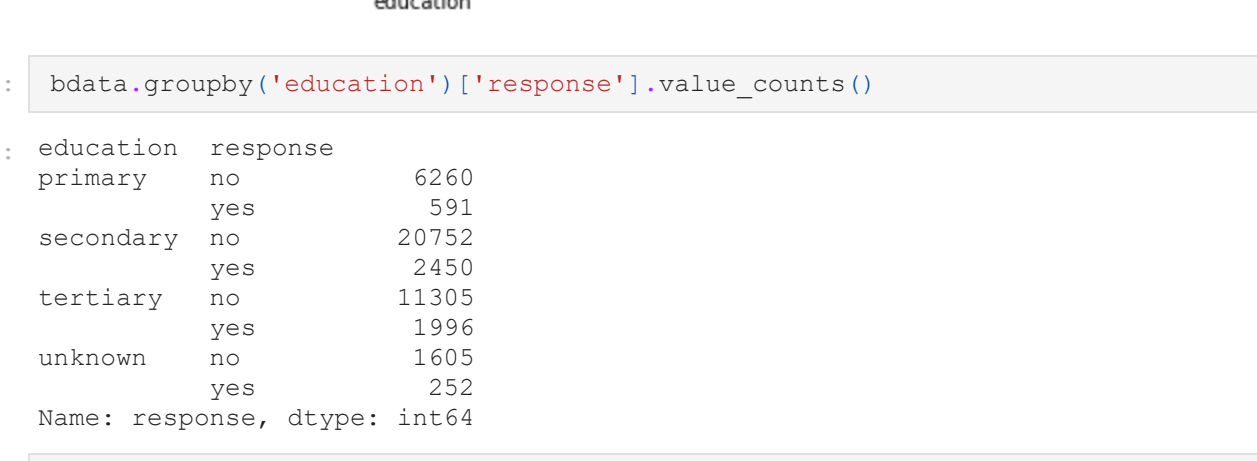
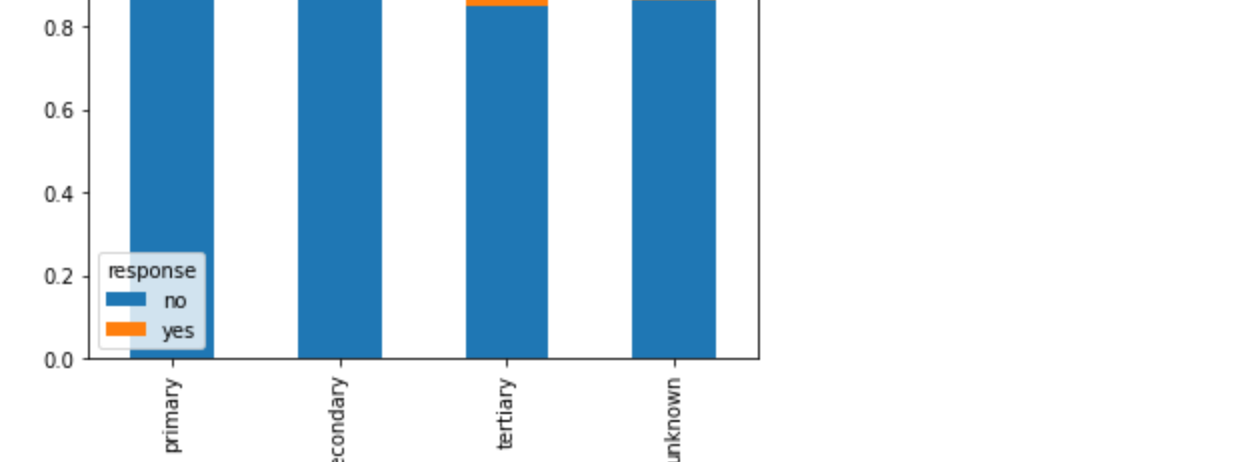
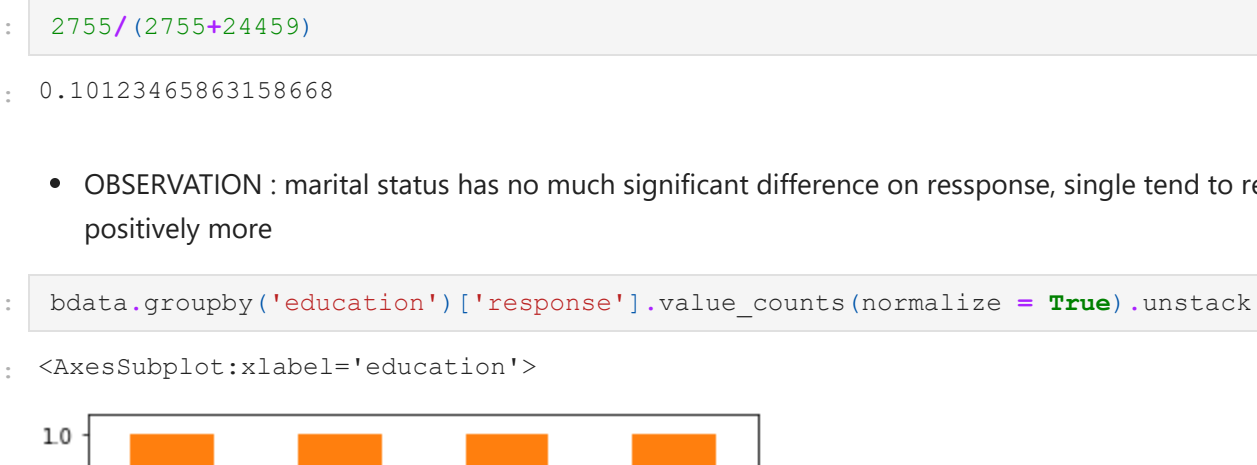
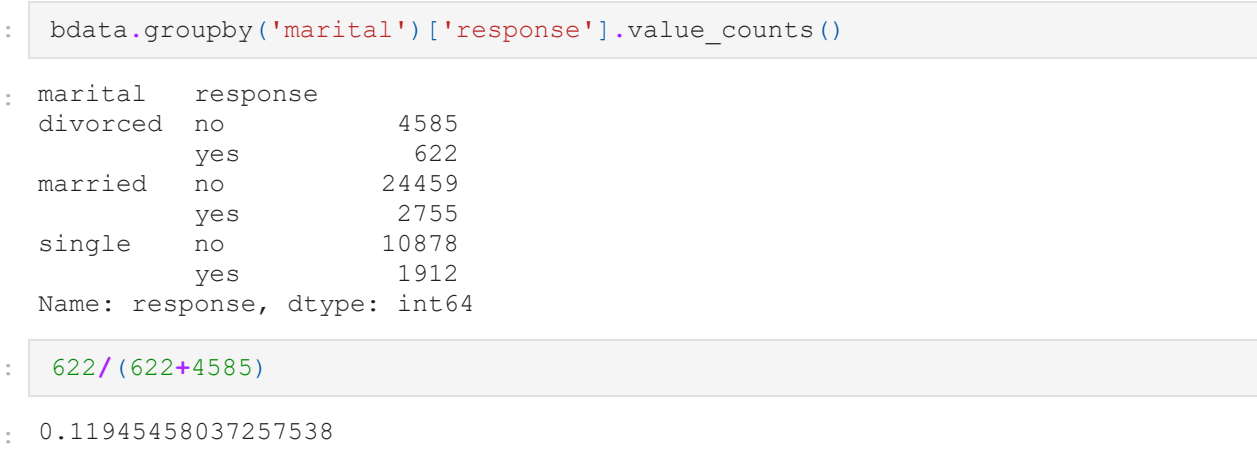
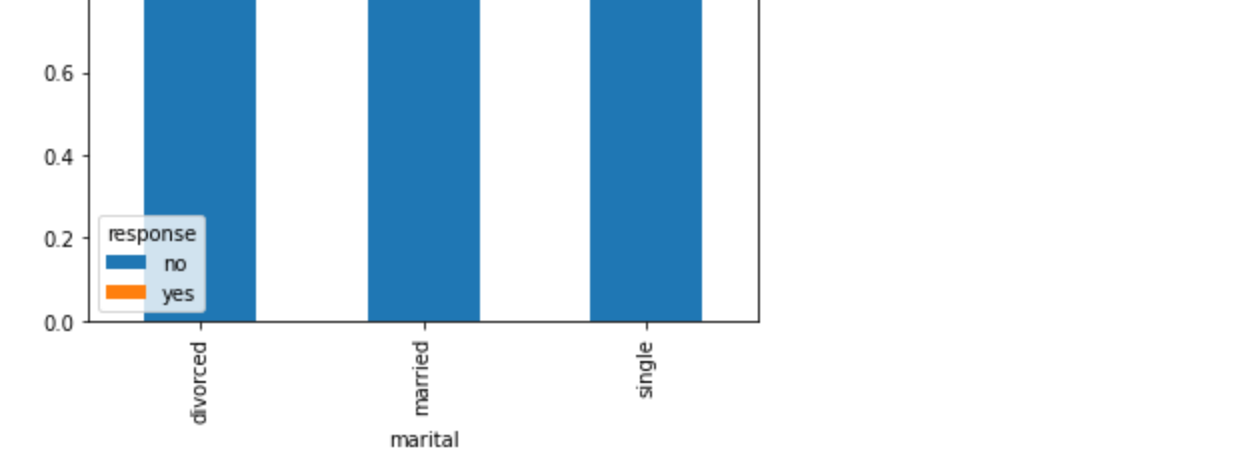


- Observation : student > retired > unemployed > management > admin

EFFECTS OF SALARY ON RESPONSE RATE



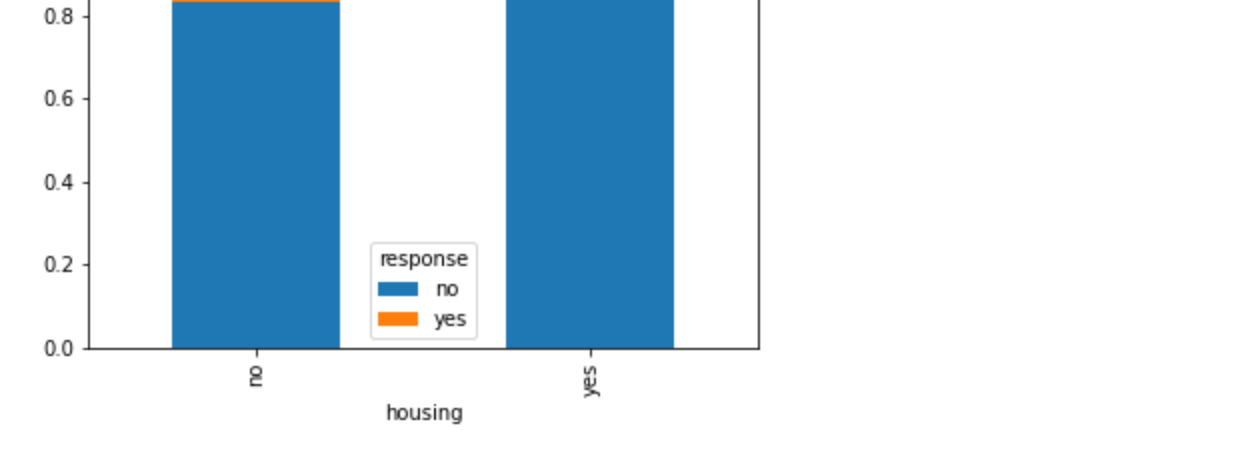
EFFECTS OF MARITAL STATUS ON RESPONSE RATE



EFFECT of account balance on response variable

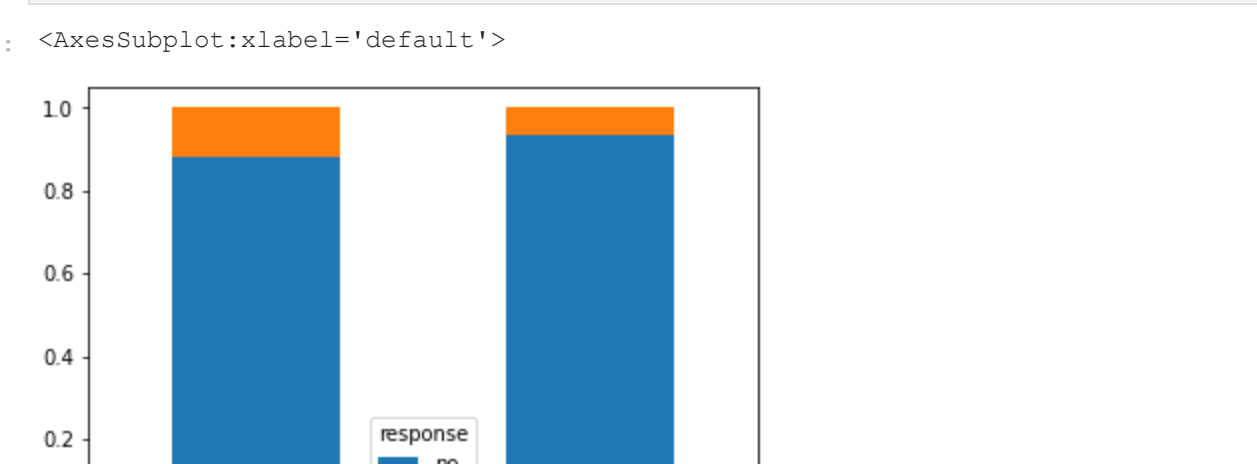
tertiary educated group has most median account balance and this group tends to more positive response than other

EFFECTS OF HOUSING LOAN ON TERM DEPOSIT

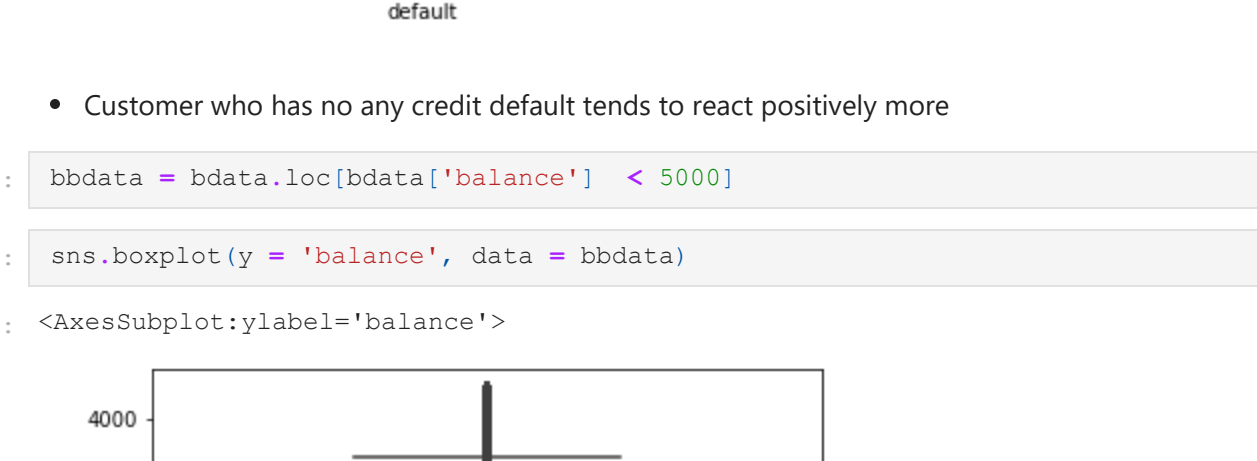


- OBSERVATION : if a customer has no housing loan tend to give more psitive response

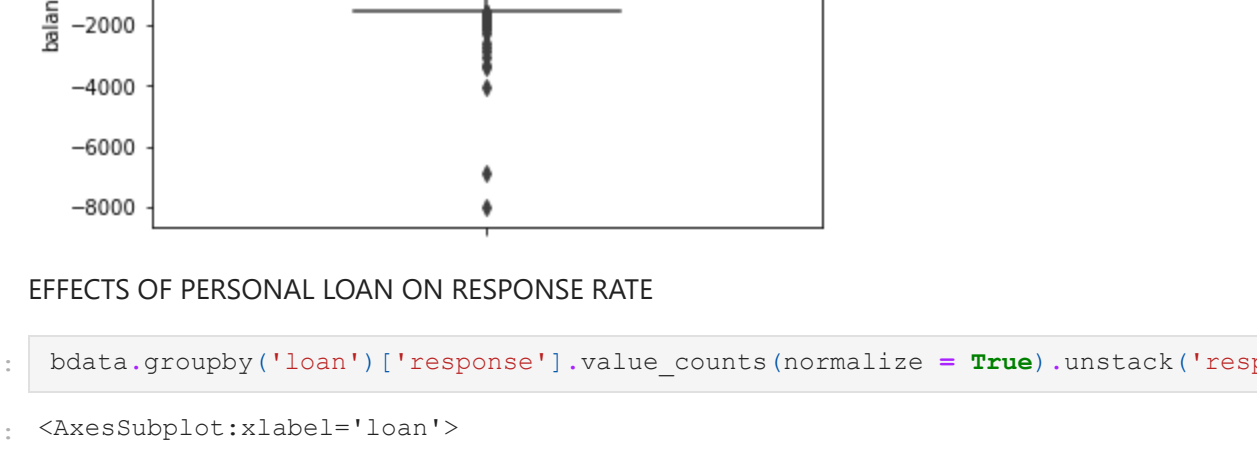
EFFECT OF CREDIT DEFAULT ON RESPONSE RATE



- Customer who has no any credit default tends to react positively more



EFFECTS OF PERSONAL LOAN ON RESPONSE RATE



- customer who has no loan tends to respond positively more than double than who has loan

current campaign effect

DOES COMMUNICATION TYPE HOLD ANY EFFECT ON RESPONSE VARIABLE

