

Data Insights

- By - Gautam Sharma

Import all required libraries

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import datetime
5 from matplotlib.ticker import FuncFormatter
6 import math as ma
7 import warnings
8 warnings.filterwarnings('ignore')
```

In [2]:

```
1 data = pd.ExcelFile('KPMG_VI_New_raw_data_update_final.xlsx')
```

In [3]:

```
1 data.sheet_names
```

Out[3]:

```
['Title Sheet',
 'Transactions',
 'NewCustomerList',
 'CustomerDemographic',
 'CustomerAddress']
```

In [4]:

```
1 df = pd.read_excel('KPMG_VI_New_raw_data_update_final.xlsx',3)
```

In [5]:

```
1 df.head()
```

Out[5]:

Note: The data and information in this document is reflective of a hypothetical situation and client. This document is to be used for KPMG Virtual Internship purposes only.

	Unnamed: 1	Unnamed: 2	Unnamed: 3		Unnamed: 4	Unnamed: 5
0	customer_id	first_name	last_name	gender	past_3_years_bike_related_purchases	DOE
1	1	Laraine	Medendorp	F	93	1953-10-12 00:00:00
2	2	Eli	Bockman	Male	81	1980-12-16 00:00:00
3	3	Arlin	Dearle	Male	61	1954-01-20 00:00:00
4	4	Talbot	NaN	Male	33	1961-10-03 00:00:00

5 rows × 26 columns



In [6]:

```
1 df.isnull().sum()
```

Out[6]:

Note: The data and information in this document is reflective of a hypothetical situation and client. This document is to be used for KPMG Virtual Internship purposes only. 0

```
Unnamed: 1
0
Unnamed: 2
125
Unnamed: 3
0
Unnamed: 4
0
Unnamed: 5
87
Unnamed: 6
506
Unnamed: 7
656
Unnamed: 8
0
Unnamed: 9
0
Unnamed: 10
302
Unnamed: 11
0
Unnamed: 12
87
Unnamed: 13
4001
Unnamed: 14
4001
Unnamed: 15
4001
Unnamed: 16
4001
Unnamed: 17
4001
Unnamed: 18
4001
Unnamed: 19
4001
Unnamed: 20
4001
Unnamed: 21
4001
Unnamed: 22
4001
Unnamed: 23
4001
Unnamed: 24
4001
Unnamed: 25
4001
dtype: int64
```

In [7]:

```
1 # rename for easier analysis
2 df.rename(columns={"Note: The data and information in this document is reflecti
3 df.rename(columns={"Unnamed: 1":"fname",
4                     "Unnamed: 2":"lname",
5                     "Unnamed: 3":"gender",
6                     "Unnamed: 4":"3y_bike_purchases",
7                     "Unnamed: 5":"DOB",
8                     "Unnamed: 6":"JT"}, inplace = True)
9 df.rename(columns={"Unnamed: 7":"Category",
10                   "Unnamed: 8":"wealth_segement",
11                   "Unnamed: 9":"D_Indicator",
12                   "Unnamed: 10":"default",
13                   "Unnamed: 11":"owns_car",
14                   "Unnamed: 12":"tencure"}, inplace = True)
15 df=df.iloc[1:]
16 df.head()
```

Out[7]:

	customer_id	fname	lname	gender	3y_bike_purchases	DOB	JT	Catego
1	1	Laraine	Medendorp	F	93	1953-10-12 00:00:00	Executive Secretary	Hea
2	2	Eli	Bockman	Male	81	1980-12-16 00:00:00	Administrative Officer	Financ Servic
3	3	Arlin	Dearle	Male	61	1954-01-20 00:00:00	Recruiting Manager	Prope
4	4	Talbot	NaN	Male	33	1961-10-03 00:00:00	NaN	
5	5	Sheila-kathryn	Calton	Female	56	1977-05-13 00:00:00	Senior Editor	Ni

5 rows × 26 columns

In [8]:

```
1 df.columns
```

Out[8]:

```
Index(['customer_id', 'fname', 'lname', 'gender', '3y_bike_purchases', 'DOB', 'JT', 'Category', 'wealth_segement', 'D_Indicator', 'default', 'owns_car', 'tencure', 'Unnamed: 13', 'Unnamed: 14', 'Unnamed: 15', 'Unnamed: 16', 'Unnamed: 17', 'Unnamed: 18', 'Unnamed: 19', 'Unnamed: 20', 'Unnamed: 21', 'Unnamed: 22', 'Unnamed: 23', 'Unnamed: 24', 'Unnamed: 25'], dtype='object')
```

In [9]:

```
1 df = df.drop(['Unnamed: 13', 'Unnamed: 14', 'Unnamed: 15', 'Unnamed: 16', 'Unnamed: 17'])
```

In [10]:

```
1 df.head()
```

Out[10]:

	customer_id	fname	lname	gender	3y_bike_purchases	DOB	JT	Category
1	1	Laraine	Medendorp	F	93	1953-10-12 00:00:00	Executive Secretary	Health
2	2	Eli	Bockman	Male	81	1980-12-16 00:00:00	Administrative Officer	Financial Services
3	3	Arlin	Dearle	Male	61	1954-01-20 00:00:00	Recruiting Manager	Professional
4	4	Talbot	NaN	Male	33	1961-10-03 00:00:00	NaN	NaN
5	5	Sheila-kathryn	Calton	Female	56	1977-05-13 00:00:00	Senior Editor	Non-Executive

In [11]:

```
1 df.shape
```

Out[11]:

(4000, 13)

In [12]:

```
1 df.isnull().sum()
```

Out[12]:

```
customer_id      0
fname            0
lname          125
gender           0
3y_bike_purchases  0
DOB              87
JT              506
Category         656
wealth_segement   0
D_Indicator       0
default          302
owns_car          0
tencure           87
dtype: int64
```

In [13]:

```
1 # Analyze the gender VS bike bought
2 gender = [0,0,0]
3 for each in df['gender']:
4     if each[0] == 'F':
5         gender[0] += 1
6     elif each[0] == "M":
7         gender[1] += 1
8     else:
9         gender[2] += 1
10
11 print(gender)
```

[2039, 1873, 88]

In [14]:

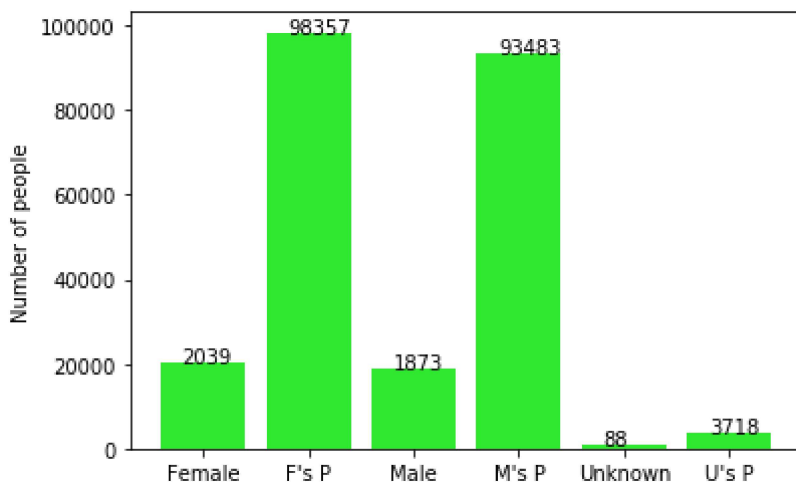
```
1 def gf(x, pos):
2     #'The two args are the value and gender'
3     return int(x)
```

In [15]:

```

1 bike = [0,0,0]
2 bike[0] += df['3y_bike_purchases'][df['gender'] == 'Female'].sum() + df['3y_bike_purchases'][df['gender'] == 'Male'].sum() + df['3y_bike_purchases'][df['gender'] == 'U'].sum()
3 bike[1] += df['3y_bike_purchases'][df['gender'] == 'Male'].sum() + df['3y_bike_purchases'][df['gender'] == 'U'].sum()
4 bike[2] += df['3y_bike_purchases'][df['gender'] == 'U'].sum()
5
6 grapho = [gender[0], bike[0], gender[1], bike[1], gender[2], bike[2]]
7 graphl = [gender[0]*10, bike[0], gender[1]*10, bike[1], gender[2]*10, bike[2]]
8 colorr = (0.1,0.9,0.1,0.9)
9 #scale up gender by 10 for easier visualization
10 formatter = FuncFormatter(gf)
11 x = np.arange(6)
12 fig, ax = plt.subplots()
13 ax.set_ylabel('Number of people')
14 ax.yaxis.set_major_formatter(formatter)
15 plt.bar(x, graphl, color = colorr)
16 for i in range(len(gender*2)):
17     plt.text(x = i-0.2, y = graphl[i]+0.1, s = grapho[i], size = 10)
18
19 plt.xticks(x, ('Female', "F's P", 'Male', "M's P", 'Unknown', "U's P"))
20 plt.show()

```



In [16]:

```

1 # Percentage of bike bought male vs female and other
2
3 avg = [0,0,0]
4 ss = sum(bike)
5
6 for i in range(len(avg)):
7     avg[i] += round(bike[i]/ss, 3)
8

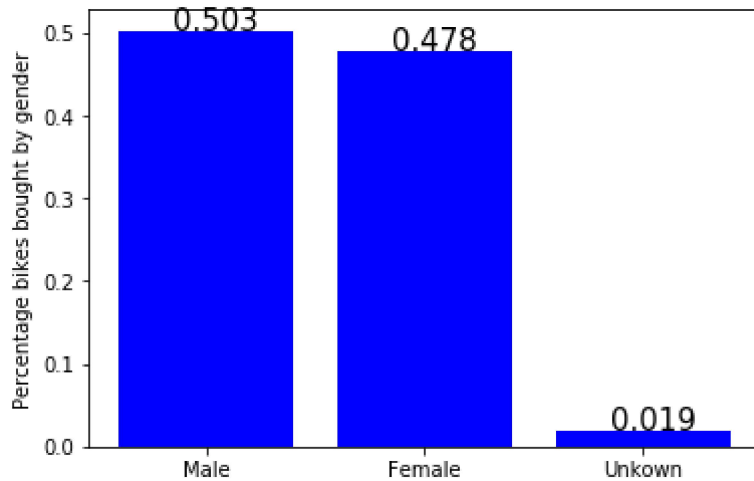
```

In [17]:

```

1 x1 = np.arange(3)
2 fig, ax1 = plt.subplots()
3 ax1.set_ylabel('Percentage bikes bought by gender')
4 plt.bar(x1, avg, color='blue')
5 for i in range(len(avg)):
6     plt.text(x = i-0.15, y=avg[i], s=avg[i], size=15)
7
8 plt.xticks(x1, ('Male', 'Female', 'Unkown'))
9 plt.show()

```



In [18]:

```

1 # Analyze on the age VS bikes bought
2 # Need to transform
3
4 print(df['DOB'][1].ctime().split(" ")[4])
5 df['Age'] = 0
6 lenn = len(df["DOB"])
7 k = 0

```

1953

In [19]:

```

1 for i in range(1, lenn):
2     if isinstance(df["DOB"][i], datetime.date):
3         t1 = len(df["DOB"][i].ctime().split(" "))
4         df['Age'][i] += int(2019-int(df['DOB'][i].ctime().split(' ')[t1-1]))
5     elif isinstance(df['DOB'][i], str):
6         t1 = len(df['DOB'][i].split("-"))
7         df['Age'][i] += int(209 - int(df['DOB'][i].split('-')[t1-1]))

```


In []:

```
1 print(k)
```

In [20]:

```
1 df.head()
```

Out[20]:

er_id	fname	Iname	gender	3y_bike_purchases	DOB	JT	Category	wealth_s
1	Laraine	Medendorp	F	93	1953-10-12 00:00:00	Executive Secretary	Health	Mass (
2	Eli	Bockman	Male	81	1980-12-16 00:00:00	Administrative Officer	Financial Services	Mass (
3	Arlin	Dearle	Male	61	1954-01-20 00:00:00	Recruiting Manager	Property	Mass (
4	Talbot	NaN	Male	33	1961-10-03 00:00:00	NaN	IT	Mass (
5	Sheila-kathryn	Calton	Female	56	1977-05-13 00:00:00	Senior Editor	NaN	Affluent (

In [21]:

```
1 #same index as above
2 ngenage = [0,0,0]
3 j       = 1
4 ss      = 0
5 stdv    = []
6 for each in df['gender']:
7     if each[0] == "F" and df['Age'][j] != 0:
8         ss += df['Age'][j]
9         stdv.append((df['Age'][j]))
10        ngenage[0] += 1
11    elif each[0] == "M" and df['Age'][j] != 0:
12        ss += df['Age'][j]
13        stdv.append((df['Age'][j]))
14        ngenage[1] += 1
15    elif df['Age'][j] != 0:
16        ss += df['Age'][j]
17        stdv.append((df['Age'][j]))
18        ngenage[2] += 1
19    j += 1
```

In [22]:

```
1 print(ss)
2 print(ngenage) # as we can see, unknown gender will unlikely to have age, don'
3 #average age not counting 0 is
4 mean_val = round(ss/sum(ngenage),0)
5 print(mean_val)
6 stdv_val = round(ma.sqrt(1/(sum(ngenage)-1)*sum((stdv - (ss/sum(ngenage)))**2)))
7 print(stdv_val)
```

164210

[2039, 1872, 1]

42.0

13.0

In [23]:

```

1  #prurchases from age 42 - 34/2, 42, 42 + 34/2
2  age_dict = {}
3  f1 = []
4  bf1 = []
5  f2 = []
6  bf2 = []
7  f3 = []
8  bf3 = []
9  f4 = []
10 bf4 = []
11 m1 = []
12 bm1 = []
13 m2 = []
14 bm2 = []
15 m3 = []
16 bm3 = []
17 m4 = []
18 bm4 = []
19 fq = mean_val - stdv_val/2
20 sq = mean_val
21 tq = mean_val + stdv_val/2
22 print(fq, sq, tq)
23 jjj = 1
24 for each in df['gender']:
25     temp = int(df['Age'][jjj])
26     bkt = int(df['3y_bike_purchases'][jjj])
27     if each[0] == "F" and temp != 0:
28         if(temp <= fq):
29             f1.append(temp)
30             bf1.append(bkt)
31         elif(fq < temp and temp <= sq):
32             f2.append(temp)
33             bf2.append(bkt)
34         elif(sq < temp and temp <= tq):
35             f3.append(temp)
36             bf3.append(bkt)
37         elif(tq < temp):
38             f4.append(temp)
39             bf4.append(bkt)
40     elif each[0] == "M" and temp != 0:
41         if(temp <= fq):
42             m1.append(temp)
43             bm1.append(bkt)
44         elif(fq < temp and temp <= sq):
45             m2.append(temp)
46             bm2.append(bkt)
47         elif(sq < temp and temp <= tq):
48             m3.append(temp)
49             bm3.append(bkt)
50         elif(tq < temp):
51             m4.append(temp)
52             bm4.append(bkt)
53     jjj += 1
54 dtt = {"Female1":f1,
55        "Female2":f2,
56        "Female3":f3,
57        "Female4":f4,
58        "Male1":m1,
59        "Male2":m2,

```

```

60         "Male3":m3,
61         "Male4":m4,
62     }

```

35.5 42.0 48.5

In [24]:

```

1 print(len(dtt["Female1"]), len(dtt["Female2"]), len(dtt["Female3"]), len(dtt["Female4"]))
2 print(len(dtt["Male1"]), len(dtt["Male2"]), len(dtt["Male3"]), len(dtt["Male4"]))

```

630 418 380 611

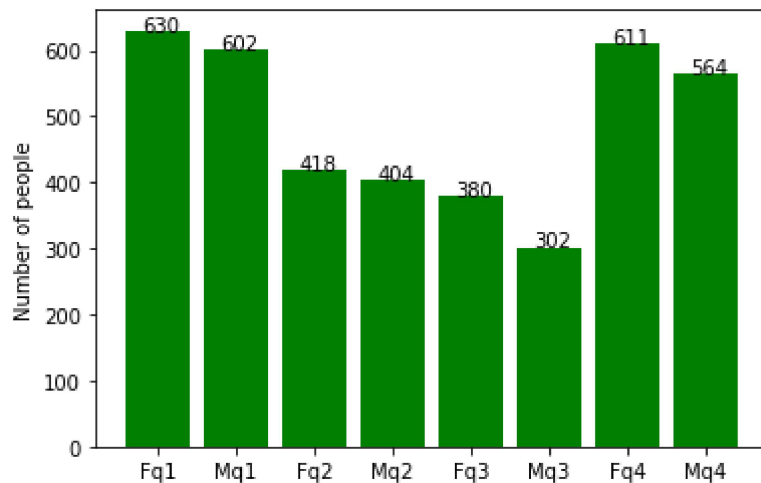
602 404 302 564

In [25]:

```

1 graphgen = [len(dtt["Female1"]), len(dtt["Male1"]), len(dtt["Female2"]), len(dtt["Male2"]), len(dtt["Female3"]), len(dtt["Male3"]), len(dtt["Female4"]), len(dtt["Male4"])]
2 colorr = (0.7,0.7,0.7,0.7)
3
4 x2 = np.arange(8)
5 fig2, ax2 = plt.subplots()
6 ax2.set_ylabel('Number of people')
7 plt.bar(x2, graphgen, color = 'green')
8 for i in range(len(graphgen)):
9     plt.text(x = i-0.2, y = graphgen[i]+0.1, s = graphgen[i], size = 10)
10
11 plt.xticks(x2, ('Fq1', 'Mq1', 'Fq2', 'Mq2', 'Fq3', 'Mq3', 'Fq4', 'Mq4'))
12 plt.show()

```



In [26]:

```
1 #check_unique() ["Category"]  
2 df["Category"].value_counts()
```

Out[26]:

```
Manufacturing      799  
Financial Services  774  
Health             602  
Retail             358  
Property           267  
IT                 223  
Entertainment      136  
Agriculture        113  
Telecommunications  72  
Name: Category, dtype: int64
```

In [27]:

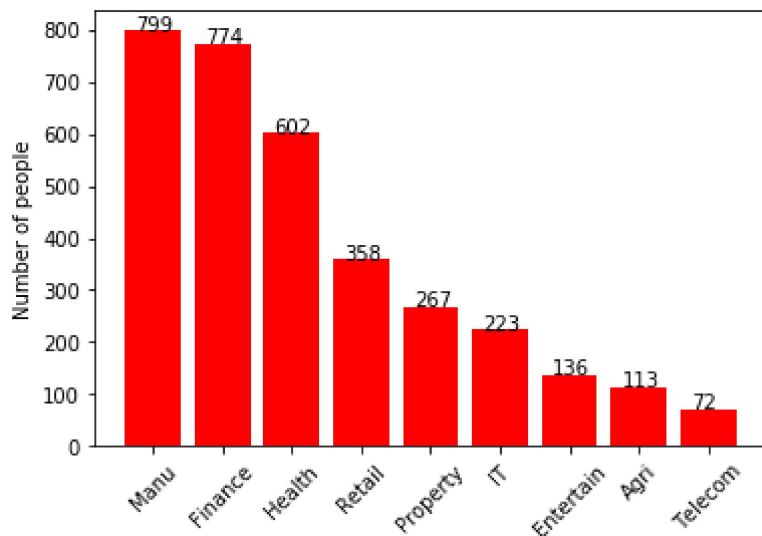
```
1 df.describe()
```

Out[27]:

	Age
count	4000.000000
mean	41.052500
std	14.103084
min	0.000000
25%	31.750000
50%	42.000000
75%	51.000000
max	188.000000

In [28]:

```
1 val = [799, 774, 602, 358, 267, 223, 136, 113, 72]
2 colorr = (0.2,0.3,0.4,0.5)
3
4 x3 = np.arange(9)
5 fig3, ax3 = plt.subplots()
6 ax3.set_ylabel('Number of people')
7 plt.bar(x3, val, color = 'red', width = 0.8)
8 for i in range(len(val)):
9     plt.text(x = i-0.25, y = val[i]+0.1, s = val[i], size = 10)
10
11 plt.xticks(x3, ("Manu", "Finance", "Health", "Retail", "Property", "IT", "Enter
12 plt.show()
```



In []:

```
1 df["wealth_segement"].value_counts()
```

In [29]:

```

1  #split into 3: M (Mass), H(High), A(Affluent)
2  jjj = 1
3  wsm = {"q1":[], "q2":[], "q3":[], "q4":[]}
4  wsh = {"q1":[], "q2":[], "q3":[], "q4":[]}
5  wsa = {"q1":[], "q2":[], "q3":[], "q4":[]}
6
7  for each in df['wealth_segement']:
8      temp = int(df['Age'][jjj])
9      if each[0] == "M" and temp != 0:
10         if(temp <= fq):
11             wsm["q1"].append(temp)
12         elif(fq < temp and temp <= sq):
13             wsm["q2"].append(temp)
14         elif(sq < temp and temp <= tq):
15             wsm["q3"].append(temp)
16         elif(tq < temp):
17             wsm["q4"].append(temp)
18     elif each[0] == "H" and temp != 0:
19         if(temp <= fq):
20             wsh["q1"].append(temp)
21         elif(fq < temp and temp <= sq):
22             wsh["q2"].append(temp)
23         elif(sq < temp and temp <= tq):
24             wsh["q3"].append(temp)
25         elif(tq < temp):
26             wsh["q4"].append(temp)
27     elif each[0] == "A" and temp != 0:
28         if(temp <= fq):
29             wsa["q1"].append(temp)
30         elif(fq < temp and temp <= sq):
31             wsa["q2"].append(temp)
32         elif(sq < temp and temp <= tq):
33             wsa["q3"].append(temp)
34         elif(tq < temp):
35             wsa["q3"].append(temp)
36     jjj += 1
37 print(len(wsm["q1"]), len(wsm["q2"]), len(wsm["q3"]), len(wsm["q4"]))
38 print(len(wsh["q1"]), len(wsh["q2"]), len(wsh["q3"]), len(wsh["q4"]))
39 print(len(wsa["q1"]), len(wsa["q2"]), len(wsa["q3"]), len(wsa["q4"]))

```

```

613 425 327 589
307 217 170 302
312 180 470 0

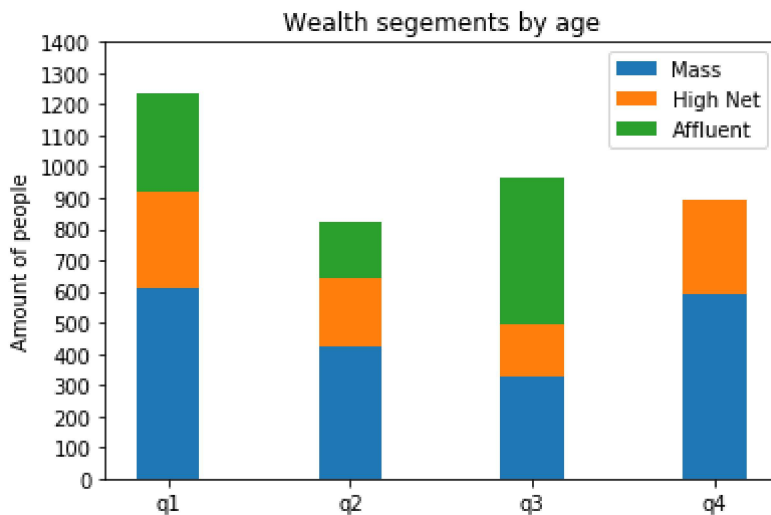
```

In [30]:

```

1 N = 4
2 wsmtp = [len(wsm["q1"]), len(wsm["q2"]), len(wsm["q3"]), len(wsm["q4"])]
3 wshtp = [len(wsh["q1"]), len(wsh["q2"]), len(wsh["q3"]), len(wsh["q4"])]
4 wsatp = [len(wsa["q1"]), len(wsa["q2"]), len(wsa["q3"]), len(wsa["q4"])]
5
6 bars = np.add(wsmtp, wshtp).tolist()
7 r = [0,1,2,3,4]
8
9 ind = np.arange(N) # the x locations for the groups
10 width = 0.35 # the width of the bars: can also be len(x) sequence
11
12 p1 = plt.bar(ind, wsmtp, width)
13 p2 = plt.bar(ind, wshtp, width, bottom=wsmtp)
14 p3 = plt.bar(ind, wsatp, width, bottom=bars)
15
16 plt.ylabel('Amount of people')
17 plt.title('Wealth segements by age')
18 plt.xticks(ind, ('q1', 'q2', 'q3', 'q4'))
19 plt.yticks(np.arange(0, 1500, 100))
20 plt.legend((p1[0], p2[0], p3[0]), ('Mass', 'High Net', 'Affluent'))
21
22 plt.show()

```



In [31]:

```
1 df['owns_car'].value_counts()
```

Out[31]:

Yes 2024

No 1976

Name: owns_car, dtype: int64

In []:

```
1
```