

```
In [1]: # Importing basic libraries

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib

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

```
In [2]: # Loading the dataset

df = pd.read_csv('https://raw.githubusercontent.com/nanthasn/Black-Friday-Sales-Prediction/master/Data/BlackFridaySales.csv')
```

```
In [3]: # top 5 rows of the df

df.head()
```

```
[3]: # top 5 rows of the df
df.head()
```

User_ID	Product_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase	
0	1000001	P0006942	F	0-17	10	A	2	0	3	NaN	NaN	8370
1	1000001	P0024842	F	0-17	10	A	2	0	1	6.0	14.0	15200
2	1000001	P00067842	F	0-17	10	A	2	0	12	NaN	NaN	1422
3	1000001	P0006542	F	0-17	10	A	2	0	12	14.0	NaN	1057
4	1000002	P0026542	M	55+	16	C	4+	0	8	NaN	NaN	7969

```
In [4]: # Getting shape of the dataset
df.shape
```

```
Out[4]: (550068, 12)
```

there are 550068 rows and 12 columns present in the dataset

```
In [5]: # Getting rows of the dataset
```

```
In [4]: # Getting shape of the dataset

df.shape
```

Out[4]: (550068, 12)

There are 550068 rows and 12 columns present in the dataset.

```
In [5]: # Getting names of the columns

df.columns
```

```
Out[5]: Index(['User_ID', 'Product_ID', 'Gender', 'Age', 'Occupation', 'City_Category',
        'Stay_In_Current_City_Years', 'Marital_Status', 'Product_Category_1',
        'Product_Category_2', 'Product_Category_3', 'Purchase'],
        dtype='object')
```

```
In [6]: # Count of types of data type

df.dtypes.value_counts()
```

```
Out[6]: int64      5
        object     5
        float64    4
        dtype: int64
```

```
In [7]: # Basic information about dataset

df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 550068 entries, 0 to 550067
Data columns (total 12 columns):
 #   Column                Non-Null Count  Dtype  ---
 0   User_ID               550068 non-null  int64  
 1   Product_ID            550068 non-null  object  
 2   Gender                550068 non-null  object  
 3   Age                   550068 non-null  object  
 4   Occupation            550068 non-null  int64  
 5   City_Category         550068 non-null  object  
 6   Stay_In_Current_City_Years  550068 non-null  object  
 7   Marital_Status        550068 non-null  int64  
 8   Product_Category_1     550068 non-null  int64  
 9   Product_Category_2     274329 non-null  float64  
10   Product_Category_3     166821 non-null  float64  
11   Purchase              550068 non-null  int64  
dtypes: float64(2), int64(5), object(5)
memory usage: 56.4+ MB
```

Age should be converted to numerical column City_Category we can convert this to a numerical column and should look at the frequency of each city category. Gender has two values and should be converted to binary values Product_Category_2 and Product_Category_3 have null values

```
In [8]: # Checking null values

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

```
Out[8]: User_ID      0
        Product_ID  0
        Gender      0
        Age         0
        Occupation  0
        Stay_In_Current_City_Years  0
        Marital_Status  0
        Product_Category_1  0
        Product_Category_2  175039
        Product_Category_3  385247
        Purchase    0
        dtype: int64
```

Product_Category_2 and Product_Category_3 has nan values.

```
In [9]: sns.heatmap(df.isnull())
plt.show()
```



Heatmap shows that the null values are present in the columns: 1. Product_Category_2 2. Product_Category_3

```
In [10]: # Total null values in dataset

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

```
Out[10]: 550885
```

```
In [11]: # Null values in percentage

df.isnull().sum()/df.shape[0]*100
```

```
Out[11]: User_ID      0.000000
        Product_ID  0.000000
        Gender      0.000000
        Age         0.000000
        Occupation  0.000000
        Stay_In_Current_City_Years  0.000000
        Marital_Status  0.000000
        Product_Category_1  0.000000
        Product_Category_2  31.566643
        Product_Category_3  69.676599
        Purchase    0.000000
        dtype: float64
```

There are 31% of the null values in the Product_Category_2 and 69% null values in the Product_Category_3

```
In [12]: # Unique elements in each attributes

df.nunique()
```

```
Out[12]: User_ID      5891
        Product_ID  3833
        Gender        2
        Age           7
        Occupation   21
        City_Category  3
        Stay_In_Current_City_Years  5
        Marital_Status  2
        Product_Category_1  20
        Product_Category_2  17
        Product_Category_3  15
        Purchase     18105
        dtype: int64
```

We can drop User_ID and Product_ID for model prediction as it has more unique values.

```
In [13]: # Checking for duplicates

df.duplicated().sum()
```

```
Out[13]: 0
```

```
In [14]: # Dropping out unnecessary columns

df.drop(['Product_ID'],inplace=True,axis=1)
```

```
In [15]: df.head()
```

```
df.drop(['Product_1D'], inplace=True, axis=1)
```

```
In [13]: df.head()
```

```
Out[13]:
```

	User_ID	Gender	Age	Occupation	City_Category	Stay_In_Current_City_Years	Marital_Status	Product_Category_1	Product_Category_2	Product_Category_3	Purchase
0	1000001	F	0-17	10	A	2	0	3	NaN	NaN	8370

```
In [16]: # Data Cleaning

df['Stay_In_Current_City_Years'] = unique()
```

```
Out[16]: array(['2', '4+', '3', '1', '0'], dtype=object)
```

```
In [17]: def cities(value):
        if 'in' in value:
            value=value.replace('+','')
            return int(value)
        else:
            return int(value)
df['Stay_In_Current_City_Years']=df['Stay_In_Current_City_Years'].apply(cities)
```

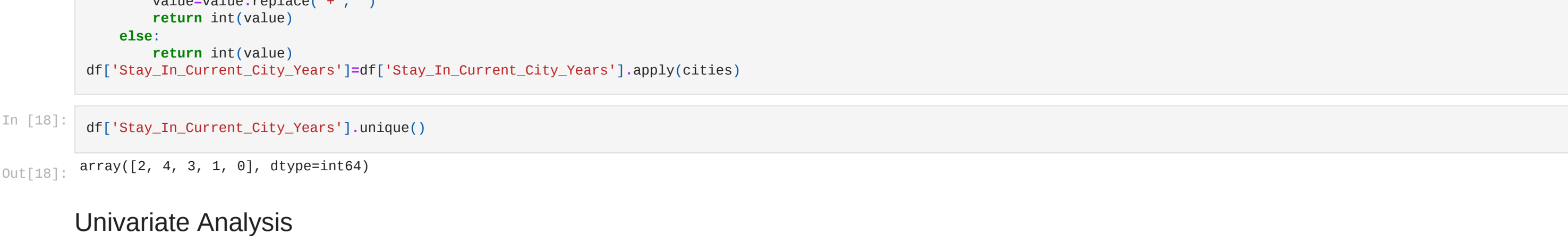
```
In [18]: df['Stay_In_Current_City_Years'].unique()
```

```
Out[18]: array([2, 4, 3, 1, 0], dtype=int64)
```

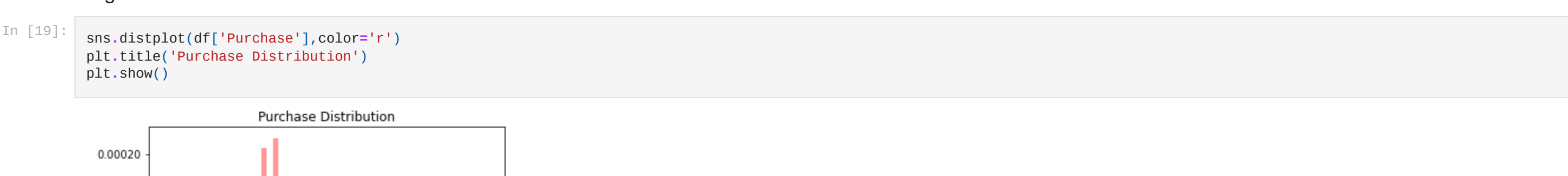
Univariate Analysis

Target Variable Purchase

```
In [19]: sns.distplot(df['Purchase'],color='r')
plt.title('Purchase Distribution')
plt.show()
```



```
In [20]: sns.boxplot(df['Purchase'])
plt.title('Boxplot of Purchase')
plt.show()
```



```
In [21]: df['Purchase'].skew()
```

```
Out[21]: 0.6001400937007128
```

```
In [22]: df['Purchase'].kurtosis()
```

```
Out[22]: -0.338377565851782
```

```
In [23]: df['Purchase'].describe()
```

```
Out[23]: count      550068.000000
        mean      9263.063713
        std       5023.065394
        min        12.000000
        25%       5823.080000
        50%       8547.080000
        75%      12054.080000
        max      23901.080000
        Name: Purchase, dtype: float64
```

Gender

```
In [24]: sns.countplot(df['Gender'])
plt.title('Countplot of Gender Feature')
plt.show()
```



```
In [25]: df['Gender'].value_counts()
```

```
Out[25]: M    414259
        F    135809
        Name: Gender, dtype: int64
```

```
In [26]: # getting percentage of the male and female

df['Gender'].value_counts(normalize=True)*100
```

```
Out[26]: M    75.318507
        F    24.681493
        Name: Gender, dtype: float64
```

There are more male than females

```
In [27]: df.groupby('Gender')['Purchase'].mean()
```

```
Out[27]: Gender
        F    8734.565765
        M    9437.526440
        Name: Purchase, dtype: float64
```

On Average the male gender spends more money on purchase.

Marital Status

```
In [28]: sns.countplot(df['Marital_Status'])
plt.show()
```



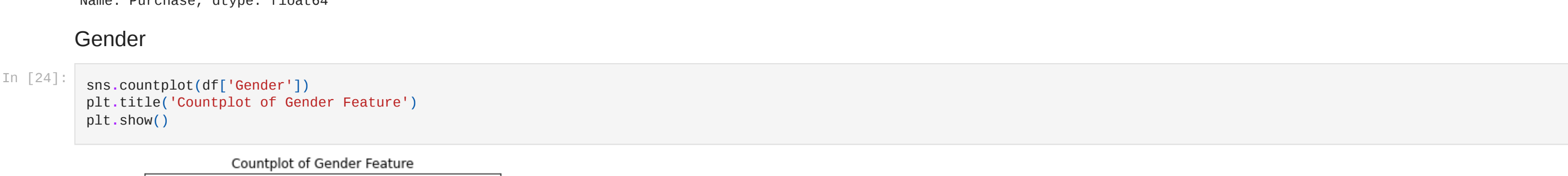
There are more unmarried people in the dataset who purchase more.

```
In [29]: df.groupby('Marital_Status')['Purchase'].mean()
```

```
Out[29]: Marital_Status
        0    9055.907419
        1    9261.174574
        Name: Purchase, dtype: float64
```

Occupation

```
In [30]: plt.figure(figsize=(18,5))
sns.countplot(df['Occupation'])
plt.show()
```



City Category

```
In [31]: sns.countplot(df['City_Category'])
plt.show()
```



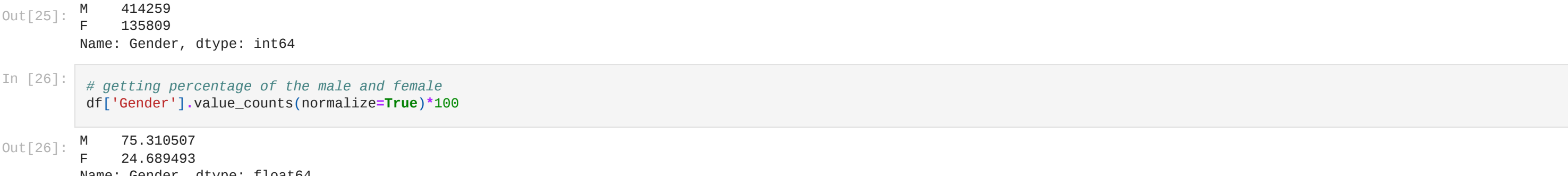
It is observed that city category B has made the most number of purchases and least is of category A.

```
In [32]: df.groupby('City_Category')['Purchase'].mean()
```

```
Out[32]: City_Category
        A    8911.939216
        B    9133.309463
        C    9719.920993
        Name: Purchase, dtype: float64
```

Stay In Current City Years

```
In [33]: sns.countplot(df['Stay_In_Current_City_Years'])
plt.show()
```

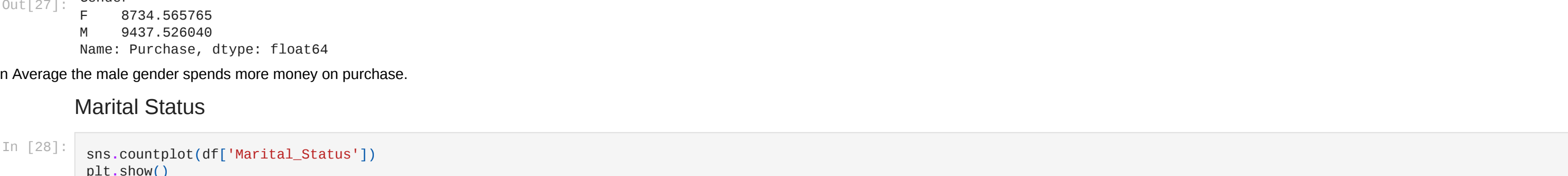


It looks like the longest someone is living in that city the less prone they are to buy new things.

```
In [34]: df.groupby('Stay_In_Current_City_Years')['Purchase'].mean()
```

```
Out[34]: Stay_In_Current_City_Years
        0    8298.075225
        1    9259.459223
        2    9129.430919
        3    9286.984119
        4    9275.595872
        Name: Purchase, dtype: float64
```

```
In [35]: df.groupby('Stay_In_Current_City_Years')['Purchase'].mean().plot(kind='bar')
```



Age

```
In [36]: df['Age'].unique()
```

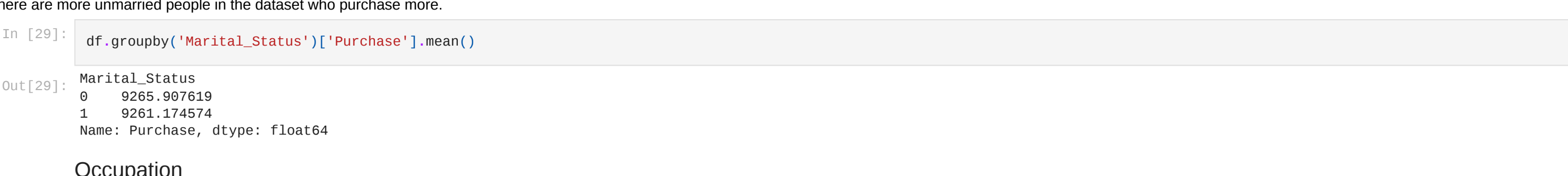
```
Out[36]: array(['0-17', '55+', '26-35', '46-50', '51-55', '36-45', '18-25'],
        dtype=object)
```

```
In [37]: def ages(value):
        if '18-25' in value:
            value=value.replace('0-17','child')
            return str(value)
        elif '36-45' in value:
            value=value.replace('26-35','adult')
            return str(value)
        elif '36-45' in value:
            value=value.replace('36-45','adult')
            return str(value)
        elif '18-25' in value:
            value=value.replace('18-25','teenage')
            return str(value)
        elif '46-50' in value:
            value=value.replace('46-50','adult')
            return str(value)
        elif '51-55' in value:
            value=value.replace('51-55','old')
            return str(value)
        else:
            value=value.replace('55+', 'old')
            return str(value)
df['Age']=df['Age'].apply(ages)
```

```
In [44]: df['Age'].unique()
```

```
Out[44]: array(['child', 'old', 'adult', 'teenage'], dtype=object)
```

```
In [46]: plt.figure(figsize=(6,6))
plt.title('Age vs purchase')
sns.barplot(x='Age',y='Purchase',data=df)
```



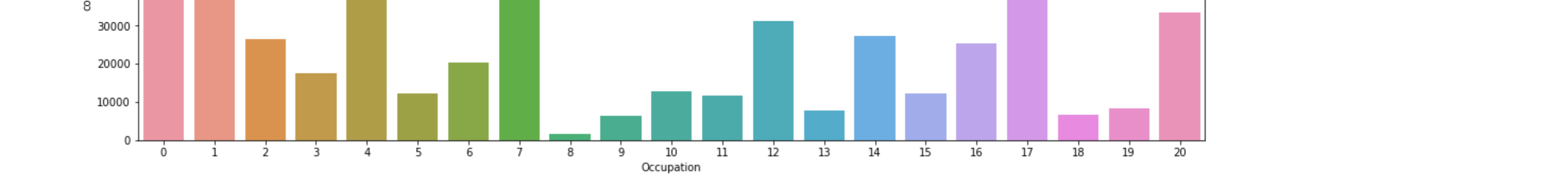
Here we understand that the purchase rate of old age is highest and lowest rate is of child.

```
In [38]: sns.countplot(df['Age'])
plt.title('Distribution of Age')
plt.xlabel('Different categories of Age')
plt.show()
```



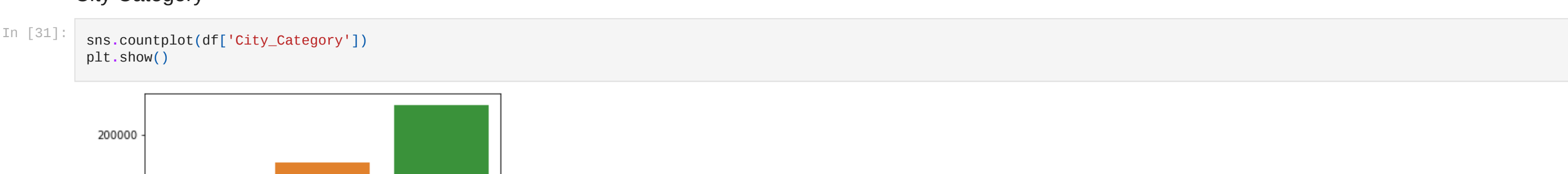
Age 26-35 group makes the most no of purchase in the group or we can say mostly adults visit the store.

```
In [39]: df.groupby('Age')['Purchase'].mean().plot(kind='bar')
```



Mean purchase rate between the age groups tends to be same except that the 51-55 age group has a little higher average purchase amount.

```
In [40]: df.groupby('Age')['Purchase'].sum().plot(kind='bar')
```



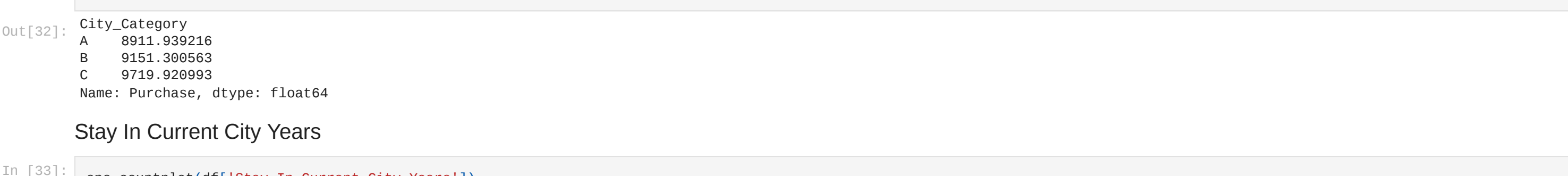
Product_Category_1

```
In [41]: plt.figure(figsize=(8,5))
sns.countplot(df['Product_Category_1'])
plt.show()
```



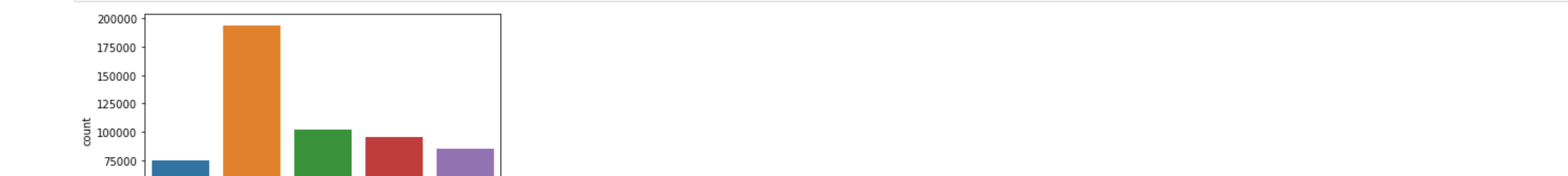
Product_Category 2

```
In [42]: plt.figure(figsize=(18,5))
sns.countplot(df['Product_Category_2'])
plt.show()
```



Product_Category 3

```
In [43]: plt.figure(figsize=(18,5))
sns.countplot(df['Product_Category_3'])
plt.show()
```



Multivariate Analysis

```
In [51]: plt.figure(figsize=(8,8))
plt.subplot(2,2,1)
plt.title('gender vs purchase')
sns.barplot(x='Gender',y='Purchase',data=df)

plt.subplot(2,2,2)
plt.title('marital_status vs purchase')
sns.barplot(x='Marital_Status',y='Purchase',hue='Gender',data=df)

plt.subplot(2,2,3)
plt.title('City vs purchase')
sns.barplot(x='City_Category',y='Purchase',hue='Gender',data=df)

plt.subplot(2,2,4)
plt.title('Age vs purchase')
sns.barplot(x='Age',y='Purchase',hue='Gender',data=df)

plt.tight_layout()
```



In [54]: plt.figure(figsize=(8,8))
sns.relplot(x='Marital_Status',y='Purchase',data=df,palette='dark',hue='Age')

```
Out[54]: <AxesSubplot: title='center: marital status vs purchase', xlabel='Marital_Status', ylabel='Purchase'>
```



1. This graph we come to know that purchase rate is more in singles rather than in married. 2. In singles purchase rate of old age is higher where as in married purchase rate of old is higher.