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

In [3]:
```python
#Read the Walmart data:
file=r'D:\walmart_data.csv'
df = pd.read_csv(file)
df
```

Out[3]:

| | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Years | Marital_Status | Product_Cate |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1000001 | P00069042 | F | 0-17 | 10 | A | 2 | 0 | |
| 1 | 1000001 | P00248942 | F | 0-17 | 10 | A | 2 | 0 | |
| 2 | 1000001 | P00087842 | F | 0-17 | 10 | A | 2 | 0 | |
| 3 | 1000001 | P00085442 | F | 0-17 | 10 | A | 2 | 0 | |
| 4 | 1000002 | P00285442 | M | 55+ | 16 | C | 4+ | 0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 550063 | 1006033 | P00372445 | M | 51-55 | 13 | B | 1 | 1 | |
| 550064 | 1006035 | P00375436 | F | 26-35 | 1 | C | 3 | 0 | |
| 550065 | 1006036 | P00375436 | F | 26-35 | 15 | B | 4+ | 1 | |
| 550066 | 1006038 | P00375436 | F | 55+ | 1 | C | 2 | 0 | |
| 550067 | 1006039 | P00371644 | F | 46-50 | 0 | B | 4+ | 1 | |

550068 rows × 10 columns

In [4]:
```python
#Checking missing values:
df.isnull().sum()/len(df)*100
```

Out[4]:
```
User_ID                       0.0
Product_ID                    0.0
Gender                        0.0
Age                           0.0
Occupation                    0.0
City_Category                 0.0
Stay_In_Current_City_Years    0.0
Marital_Status                0.0
Product_Category              0.0
Purchase                      0.0
dtype: float64
```

In [5]: 
```
#Checking the characteristics of the data:
df.describe(include='all')
```

Out[5]:

| | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Years | Marital_Status | P |
|---|---|---|---|---|---|---|---|---|---|
| count | 5.500680e+05 | 550068 | 550068 | 550068 | 550068.000000 | 550068 | 550068 | 550068.000000 | |
| unique | NaN | 3631 | 2 | 7 | NaN | 3 | 5 | NaN | |
| top | NaN | P00265242 | M | 26-35 | NaN | B | 1 | NaN | |
| freq | NaN | 1880 | 414259 | 219587 | NaN | 231173 | 193821 | NaN | |
| mean | 1.003029e+06 | NaN | NaN | NaN | 8.076707 | NaN | NaN | 0.409653 | |
| std | 1.727592e+03 | NaN | NaN | NaN | 6.522660 | NaN | NaN | 0.491770 | |
| min | 1.000001e+06 | NaN | NaN | NaN | 0.000000 | NaN | NaN | 0.000000 | |
| 25% | 1.001516e+06 | NaN | NaN | NaN | 2.000000 | NaN | NaN | 0.000000 | |
| 50% | 1.003077e+06 | NaN | NaN | NaN | 7.000000 | NaN | NaN | 0.000000 | |
| 75% | 1.004478e+06 | NaN | NaN | NaN | 14.000000 | NaN | NaN | 1.000000 | |
| max | 1.006040e+06 | NaN | NaN | NaN | 20.000000 | NaN | NaN | 1.000000 | |

In [6]: 
```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 550068 entries, 0 to 550067
Data columns (total 10 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            550068 non-null  int64
 9   Purchase                    550068 non-null  int64
dtypes: int64(5), object(5)
memory usage: 42.0+ MB
```

Initial Observations:

1. There are no missing values in the data.
2. There are 3631 unique product IDs in the dataset. P00265242 is the most sold Product ID.
3. There are 7 unique age groups and most of the purchase belongs to age 26-35 group.
4. There are 3 unique citi categories with category B being the highest.
5. 5 unique values for Stay_in_current_citi_years with 1 being the highest.
6. The difference between mean and median seems to be significant for purchase that suggests outliers in the data.
7. Minimum & Maximum purchase is 12 and 23961 suggests the purchasing behaviour is quite spread over a aignificant range of values. Mean is 9264 and 75% of purchase is of less than or equal to 12054. It suggest most of the purchase is not more than 12k.
8. Few categorical variable are of integer data type. It can be converted to character type.
9. Out of 550068 data points, 414259's gender is Male and rest are the female. Male purchase count is much higher than female.
10. Standard deviation for purchase have significant value which suggests data is more spread out for this attribute.#

In [7]: 
```
columns=['User_ID','Occupation', 'Marital_Status', 'Product_Category']
df[columns]=df[columns].astype('object')
```

In [8]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 550068 entries, 0 to 550067
Data columns (total 10 columns):
 #   Column                      Non-Null Count   Dtype
---  ------                      --------------   -----
 0   User_ID                     550068 non-null  object
 1   Product_ID                  550068 non-null  object
 2   Gender                      550068 non-null  object
 3   Age                         550068 non-null  object
 4   Occupation                  550068 non-null  object
 5   City_Category               550068 non-null  object
 6   Stay_In_Current_City_Years  550068 non-null  object
 7   Marital_Status              550068 non-null  object
 8   Product_Category            550068 non-null  object
 9   Purchase                    550068 non-null  int64
dtypes: int64(1), object(9)
memory usage: 42.0+ MB
```

In [9]: `df.describe(include='all')`

Out[9]:

| | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Years | Marital_Status | Product_ |
|---|---|---|---|---|---|---|---|---|---|
| count | 550068.0 | 550068 | 550068 | 550068 | 550068.0 | 550068 | 550068 | 550068.0 | |
| unique | 5891.0 | 3631 | 2 | 7 | 21.0 | 3 | 5 | 2.0 | |
| top | 1001680.0 | P00265242 | M | 26-35 | 4.0 | B | 1 | 0.0 | |
| freq | 1026.0 | 1880 | 414259 | 219587 | 72308.0 | 231173 | 193821 | 324731.0 | |
| mean | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| std | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| min | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 25% | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 50% | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| 75% | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |
| max | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | |

## Observation post modifying the categorical variable's data type:

1. There are 5891 unique users, and userid 1001680 being with the highest count.
2. The customers belongs to 21 distinct occupation for the purchases being made with Occupation 4 being the highest.
3. Marital status unmarried contribute more in terms of the count for the purchase.
4. There are 20 unique product categories with 5 being the highest.

In [10]:
```python
# Checking how categorical variables contributes to the entire data
categ_cols = ['Gender', 'Age', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status']
df[categ_cols].melt().groupby(['variable', 'value'])[['value']].count()/len(df)
```

Out[10]:

|  |  | value |
|---|---|---|
| **variable** | **value** | |
| **Age** | **0-17** | 0.027455 |
| | **18-25** | 0.181178 |
| | **26-35** | 0.399200 |
| | **36-45** | 0.199999 |
| | **46-50** | 0.083082 |
| | **51-55** | 0.069993 |
| | **55+** | 0.039093 |
| **City_Category** | **A** | 0.268549 |
| | **B** | 0.420263 |
| | **C** | 0.311189 |
| **Gender** | **F** | 0.246895 |
| | **M** | 0.753105 |
| **Marital_Status** | **0** | 0.590347 |
| | **1** | 0.409653 |
| **Stay_In_Current_City_Years** | **0** | 0.135252 |
| | **1** | 0.352358 |
| | **2** | 0.185137 |
| | **3** | 0.173224 |
| | **4+** | 0.154028 |

## Observations:

1. 40% of the purchase done by aged 26-35 and 78% purchase are done by the customers aged betwe en the age 18-45 (40%: 26-35, 18%: 18-25, 20%: 36-45)
2. 75% of the purchase count are done by Male and 25% by Female
3. 60% Single, 40% Married contributes to the purchase count.
4. 35% Staying in the city from 1 year, 18% from 2 years, 17% from 3 years
5. There are 20 product categories in total.
6. There are 20 different types of occupations in the city.

In [11]:
```python
#Checking how the data is spread basis distinct users

df2=df.groupby(['User_ID'])['Age'].unique()
df2.value_counts()/len(df2)
```

Out[11]:
```
Age
[26-35]    0.348498
[36-45]    0.198099
[18-25]    0.181463
[46-50]    0.090137
[51-55]    0.081650
[55+]      0.063147
[0-17]     0.037006
Name: count, dtype: float64
```

## Observation:

1. We can see 35% of the users are aged 26-35. 73% of users are aged between 18-45.
2. From the previous observation we saw 40% of the purchase are done by users aged 26-35. And, we have 35% of users aged between 26-35 and they are contributing 40% of total purchase count.S o, we can infer users aged 26-35 are more frequent  customers.

In [12]:
```python
df2=df.groupby(['User_ID'])['Gender'].unique()
df2.value_counts()/len(df2)
```

Out[12]:
```
Gender
[M]    0.717196
[F]    0.282804
Name: count, dtype: float64
```

## Observation:

1. We have 72% male users and 28% female users. Combining with previous observations we can see 72% of male users contributing to 75% of the purchase count and 28% of female users are contributing to 25% of the purchase count.

In [13]:
```python
df2=df.groupby(['User_ID'])['Marital_Status'].unique()
df2.value_counts()/len(df2)
```

Out[13]:
```
Marital_Status
[0]    0.580037
[1]    0.419963
Name: count, dtype: float64
```

## Observation:

1. We have 58% of the single users and 42% of married users. Combining with previous observation, single users contributes more as 58% of the single contributes to the 60% of the purchase count.

In [14]:
```python
df2=df.groupby(['User_ID'])['City_Category'].unique()
df2.value_counts()/len(df2)
```

Out[14]:
```
City_Category
[C]    0.532847
[B]    0.289764
[A]    0.177389
Name: count, dtype: float64
```

## Observation:

1. 53% of the users belong to city category C whereas 29% to category B and 18% belong to category A. Combining from the previous observation category B purchase count is 42% and Category C purchase count is 31%. We can clearly see category B are more actively purchasing inspite of the fact they are only 28% of the total users. On the other hand, we have 53% of category C users but they only contribute 31% of the total purchase count.

In [15]:
```python
#Checking the age group distribution in different city categories
pd.crosstab(index=df["City_Category"],columns=df["Age"],margins=True,normalize="index")
```

Out[15]:

| Age | 0-17 | 18-25 | 26-35 | 36-45 | 46-50 | 51-55 | 55+ |
|---|---|---|---|---|---|---|---|
| City_Category | | | | | | | |
| A | 0.017222 | 0.186400 | 0.499222 | 0.180185 | 0.051496 | 0.041288 | 0.024188 |
| B | 0.023511 | 0.187076 | 0.396171 | 0.205898 | 0.088272 | 0.076743 | 0.022330 |
| C | 0.041612 | 0.168705 | 0.316974 | 0.209131 | 0.103333 | 0.085649 | 0.074596 |
| All | 0.027455 | 0.181178 | 0.399200 | 0.199999 | 0.083082 | 0.069993 | 0.039093 |

## Observation:

1. We have seen earlier that city category B and A constitutes less percentage of total populat
~~ion, but they contribute more towards purchase count. We can see from above results large perce~~

```
In [16]: #Checking how genders are contributing towards toatl purchase amount
         df2=pd.DataFrame(df.groupby(['Gender'])[['Purchase']].sum())

         df2['percent'] = (df2['Purchase'] /
                           df2['Purchase'].sum()) * 100
         df2
```

Out[16]:

| Gender | Purchase | percent |
|---|---|---|
| F | 1186232642 | 23.278576 |
| M | 3909580100 | 76.721424 |

## Observation:

1. We can see male(72% of the population) contributes to more than 76% of the total purchase amount whereas female(28% of the population) contributes 23% of the total purchase amount.

```
In [17]: #Checking how purchase value are spread among differnt age categories
         df2=pd.DataFrame(df.groupby(['Age'])[['Purchase']].sum())

         df2['percent'] = (df2['Purchase'] /
                           df2['Purchase'].sum()) * 100
         df2
```

Out[17]:

| Age | Purchase | percent |
|---|---|---|
| 0-17 | 134913183 | 2.647530 |
| 18-25 | 913848675 | 17.933325 |
| 26-35 | 2031770578 | 39.871374 |
| 36-45 | 1026569884 | 20.145361 |
| 46-50 | 420843403 | 8.258612 |
| 51-55 | 367099644 | 7.203947 |
| 55+ | 200767375 | 3.939850 |

## Observation:

1. We can see the net purchase amount spread is similar to the purchase count spread among the different age groups.

```
In [18]: df2=pd.DataFrame(df.groupby(['Marital_Status'])['Purchase'].sum())

         df2['percent'] = (df2['Purchase'] /
                           df2['Purchase'].sum()) * 100
         df2
```

Out[18]:

| Marital_Status | Purchase | percent |
|---|---|---|
| 0 | 3008927447 | 59.047057 |
| 1 | 2086885295 | 40.952943 |

## Observations:

1. Single users are contributing 59% towards the total purchase amount in comparison to 41% by married users.

In [19]:
```python
df2=pd.DataFrame(df.groupby(['City_Category'])['Purchase'].sum())

df2['percent'] = (df2['Purchase'] /
                  df2['Purchase'].sum()) * 100
df2
```

Out[19]:

| City_Category | Purchase | percent |
|---|---|---|
| A | 1316471661 | 25.834381 |
| B | 2115533605 | 41.515136 |
| C | 1663807476 | 32.650483 |

## Observations:

1. City_category contribution to the total purchase amount is also similar to their contributio n towards Purchase count. Still, combining with previous observation we can City_category C alt hough has percentage purchase count of 31% but they contribute more in terms of purchase amount i.e. 32.65%. We can infer City category C purchase higher value products.

In [20]:
```python
# Users with highest number of purchases
df.groupby(['User_ID'])['Purchase'].count().nlargest(10)
```

Out[20]:
```
User_ID
1001680    1026
1004277     979
1001941     898
1001181     862
1000889     823
1003618     767
1001150     752
1001015     740
1005795     729
1005831     727
Name: Purchase, dtype: int64
```

In [21]:
```python
#Users with highest purchases amount
df.groupby(['User_ID'])['Purchase'].sum().nlargest(10)
```

Out[21]:
```
User_ID
1004277    10536909
1001680     8699596
1002909     7577756
1001941     6817493
1000424     6573609
1004448     6566245
1005831     6512433
1001015     6511314
1003391     6477160
1001181     6387961
Name: Purchase, dtype: int64
```

## Observation:

1. The users with high number of purchases contribute more to the purchase amount. Still, we ca n see there are few users not in the list of top 10 purchase counts are there in list of top 10 purchase amount. Also, the user 1004277 with lesser purchase count(979) has a much higher purch ase amount than the user(1001680) with top purchase count.

In [22]:
```python
df2=pd.DataFrame(df.groupby(['Occupation'])[['Purchase']].sum())

df2['percent'] = (df2['Purchase'] /
                  df2['Purchase'].sum()) * 100
df2
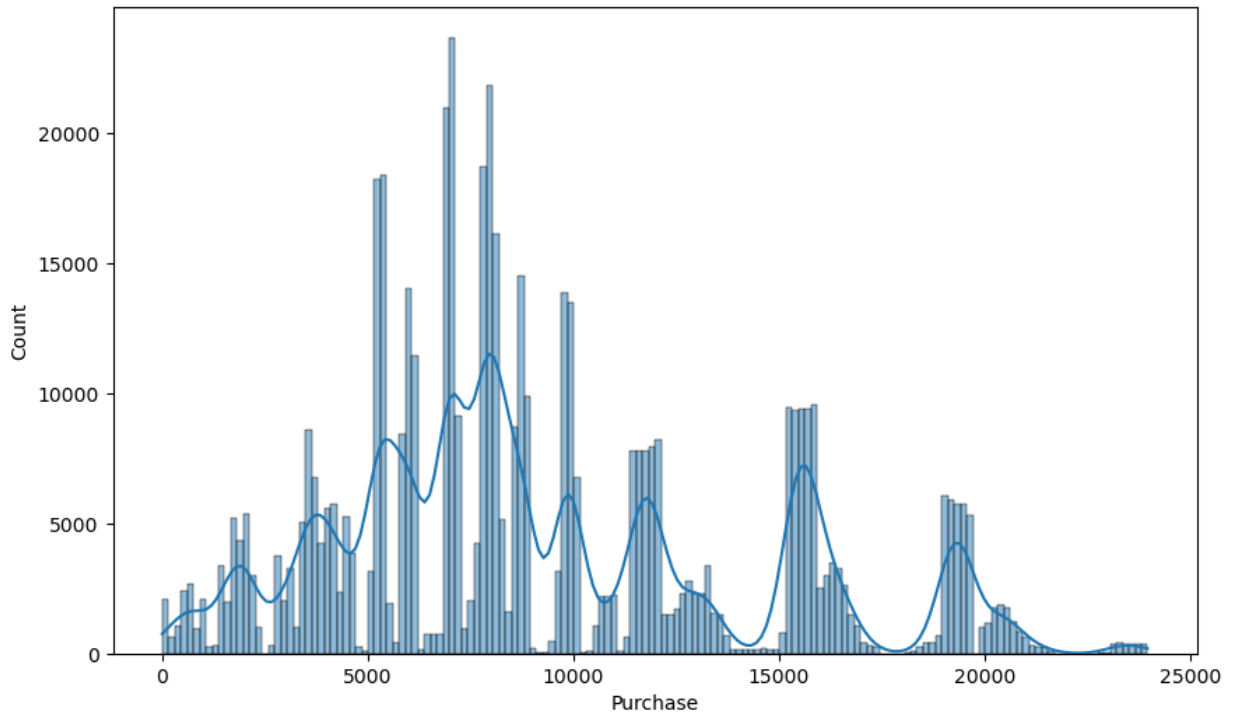```

Out[22]:

| Occupation | Purchase | percent |
|---|---|---|
| 0 | 635406958 | 12.469198 |
| 1 | 424614144 | 8.332609 |
| 2 | 238028583 | 4.671062 |
| 3 | 162002168 | 3.179123 |
| 4 | 666244484 | 13.074352 |
| 5 | 113649759 | 2.230258 |
| 6 | 188416784 | 3.697482 |
| 7 | 557371587 | 10.937835 |
| 8 | 14737388 | 0.289206 |
| 9 | 54340046 | 1.066367 |
| 10 | 115844465 | 2.273327 |
| 11 | 106751618 | 2.094889 |
| 12 | 305449446 | 5.994126 |
| 13 | 71919481 | 1.411345 |
| 14 | 259454692 | 5.091527 |
| 15 | 118960211 | 2.334470 |
| 16 | 238346955 | 4.677310 |
| 17 | 393281453 | 7.717738 |
| 18 | 60721461 | 1.191595 |
| 19 | 73700617 | 1.446298 |
| 20 | 296570442 | 5.819885 |

# # Observations:

1. Some of the Occupation like 0, 4, 7 has contributed more towards total purchase amount.

In [24]:
```python
df2=pd.DataFrame(df.groupby(['Product_Category'])[['Purchase']].sum())

df2['percent'] = (df2['Purchase'] /
                  df2['Purchase'].sum()) * 100
df2
```

Out[24]:

| Product_Category | Purchase | percent |
|---|---|---|
| 1 | 1910013754 | 37.482024 |
| 2 | 268516186 | 5.269350 |
| 3 | 204084713 | 4.004949 |
| 4 | 27380488 | 0.537313 |
| 5 | 941835229 | 18.482532 |
| 6 | 324150302 | 6.361111 |
| 7 | 60896731 | 1.195035 |
| 8 | 854318799 | 16.765114 |
| 9 | 6370324 | 0.125011 |
| 10 | 100837301 | 1.978827 |
| 11 | 113791115 | 2.233032 |
| 12 | 5331844 | 0.104632 |
| 13 | 4008601 | 0.078665 |
| 14 | 20014696 | 0.392767 |
| 15 | 92969042 | 1.824420 |
| 16 | 145120612 | 2.847840 |
| 17 | 5878699 | 0.115363 |
| 18 | 9290201 | 0.182310 |
| 19 | 59378 | 0.001165 |
| 20 | 944727 | 0.018539 |

## Observations:

1. 1, 8, 5 are among the highest yielding product categories and 19, 20, 13 are among the lowest in terms of their contribution to total amount.

In [25]:
```python
df2=pd.DataFrame(df.groupby(['Stay_In_Current_City_Years'])[['Purchase']].sum())

df2['percent'] = (df2['Purchase'] /
                  df2['Purchase'].sum()) * 100
df2
```

Out[25]:

| Stay_In_Current_City_Years | Purchase | percent |
|---|---|---|
| 0 | 682979229 | 13.402754 |
| 1 | 1792872533 | 35.183250 |
| 2 | 949173931 | 18.626547 |
| 3 | 884902659 | 17.365290 |
| 4+ | 785884390 | 15.422160 |

## Univariate Analysis:

We can explore the distribution of the data for the quantitative attributes using histplot.

In [26]:
```python
plt.figure(figsize=(10, 6))
sns.histplot(data=df, x="Purchase", kde=True)
plt.show()
```



## Observation:

1. We can see purchase value between 5000 and 10000 have higher count. From the initial observa
tion we have already seen the mean and median is 9263 and 8047 respectively. Also, we can see t
here are outliers in the data.

In [27]:
```python
plt.figure(figsize=(5, 4))
sns.boxplot(data=df, y='Purchase')
plt.show()
```
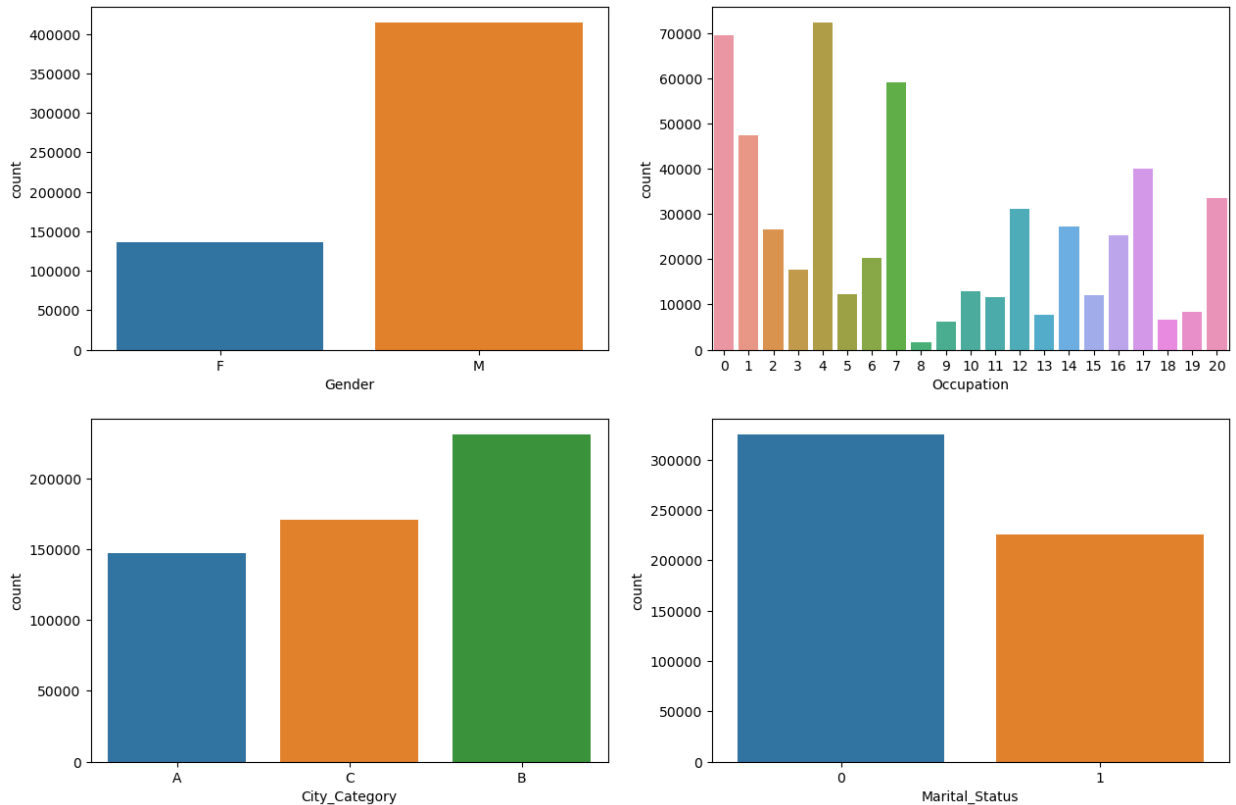


## Observation:

We can see there are outliers in the data for purchase.

Univariate analysis for qualitative variables:

In [28]:
```python
fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(15, 10))
sns.countplot(data=df, x='Gender', ax=axs[0,0])
sns.countplot(data=df, x='Occupation', ax=axs[0,1])
sns.countplot(data=df, x='City_Category', ax=axs[1,0])
sns.countplot(data=df, x='Marital_Status', ax=axs[1,1])
plt.show()
```

# Observations:

1. We can clearly see from the graphs above the purchases done by males are much higher than females.
2. We have 21 occupations categories. Occupation category 4, 0, and 7 are with higher number of purchases and category 8 with the lowest number of purchaes.
3. The purchases are highest from City category B.
4. Single customer purchases are higher than married users.

In [29]:
```python
plt.figure(figsize=(12, 5))
sns.countplot(data=df, x='Product_Category')
plt.show()
```

## Observations:

1. There are 20 product categories with product category 1, 5 and 8 having higher purchasing frequency.

In [32]: `#Bivariate Analysis:`

In [31]:
```python
fig, axs = plt.subplots(nrows=1, ncols=2, figsize=(16,5))
sns.histplot(data=df[df['Gender']=='M']['Purchase'], ax=axs[0]).set_title("Male Spending ")
sns.histplot(data=df[df['Gender']=='F']['Purchase'], ax=axs[1]).set_title("Female Spending")
plt.show()
```
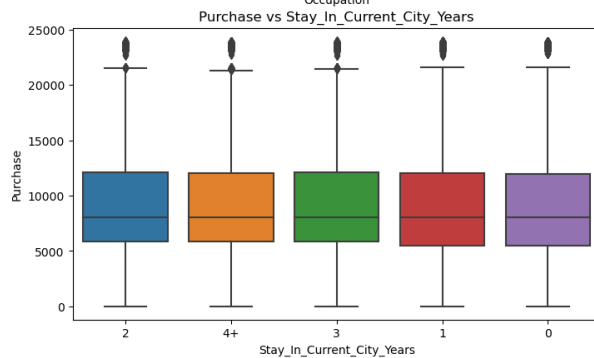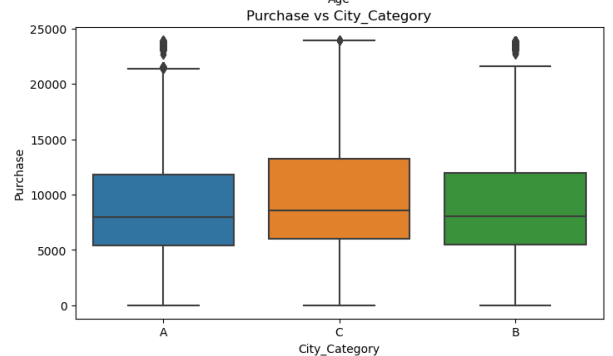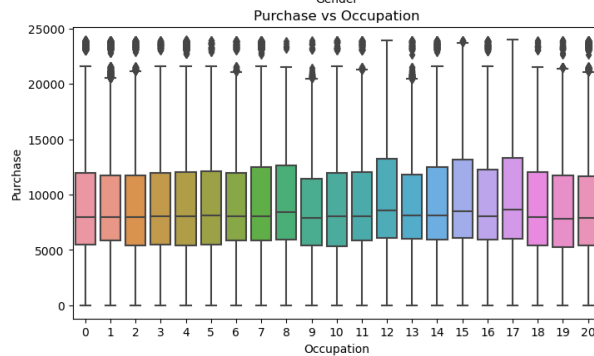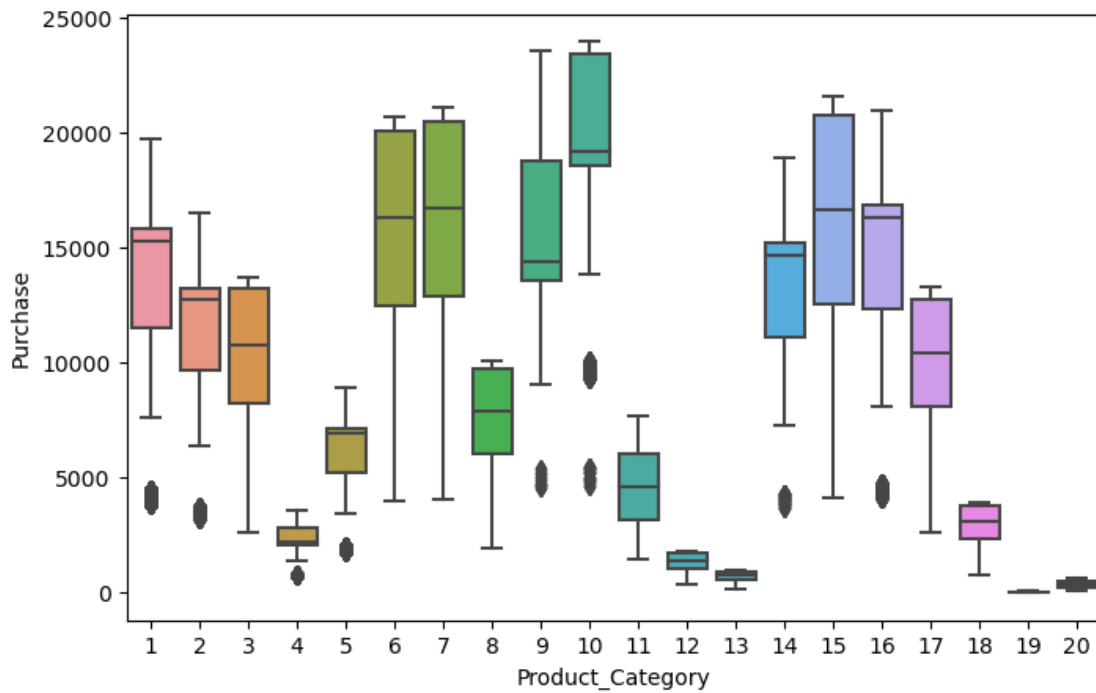


## Observations:

1. From the above histplot, we can clearly see spending behaviour is very much similar in nature for both males and females as the maximum purchase count are between the purchase value range of 5000-10000 for both. But, the purchase count are more in case of males.

In [33]:
```python
attr = ['Gender', 'Age', 'Occupation', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status']

fig, axs = plt.subplots(nrows=3, ncols=2, figsize=(18, 10))
fig.subplots_adjust(top=1.3)
count = 0
for row in range(3):
    for col in range(2):
        sns.boxplot(data=df, y='Purchase', x=attr[count], ax=axs[row, col],)
        axs[row,col].set_title(f"Purchase vs {attr[count]}")
        count += 1
plt.show()

plt.figure(figsize=(8, 5))
sns.boxplot(data=df, y='Purchase', x='Product_Category')
plt.show()
```
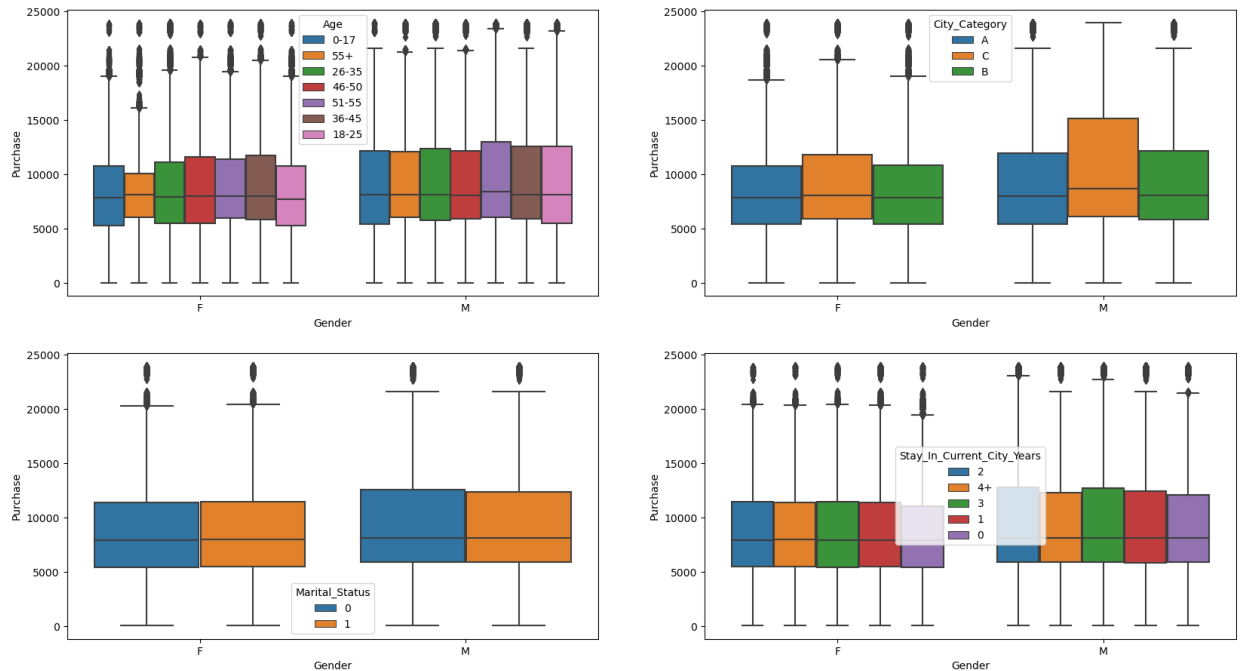
## Observations:

1. The spending behaviour for males and females are similar as we had seen from the above histp lot. Males purchasing value are in the little higher range than females.

2. Among differnt age categories, we see similar purchase behaviour. For all age groups, most o f the purchases are of the values between 5k to 12k with all have some outliers.

3. Among different occupation as well, we see similar purchasing behaviour in terms of the purc hase values.

4. Similarly for City category, stay in current city years, marital status - we see the users s pends mostly in the range of 5k to 12k.

5. We see variations among product categories. Product category 10 products are the costliest o nes. Also, there are few outliers for some of the product categories.

In [34]: *#Multivariate analysis:*

```
In [35]: fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(20, 6))
         fig.subplots_adjust(top=1.5)
         sns.boxplot(data=df, y='Purchase', x='Gender', hue='Age', ax=axs[0,0])
         sns.boxplot(data=df, y='Purchase', x='Gender', hue='City_Category', ax=axs[0,1])

         sns.boxplot(data=df, y='Purchase', x='Gender', hue='Marital_Status', ax=axs[1,0])
         sns.boxplot(data=df, y='Purchase', x='Gender', hue='Stay_In_Current_City_Years', ax=axs[1,1])

         plt.show()
```
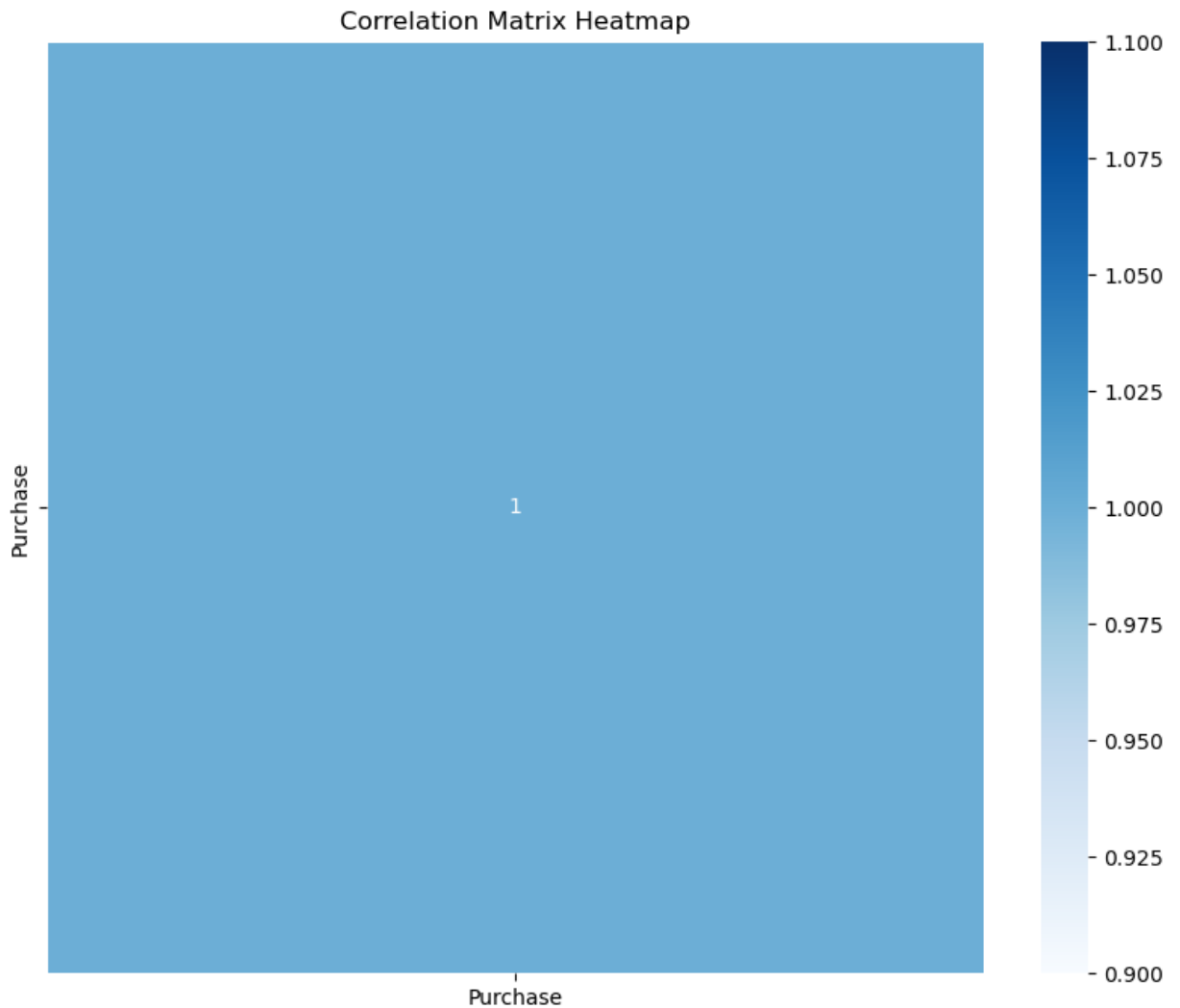


# Observations:

1. The purchasing pattern is very much similar for males and females even among differnt age groups.
2. The purchasing behaviour of males and females basis different citi categories is also similar in nature. Still, males from city category B tends to purchase costlier products in comparison to females.
3. Males and females spending behaviour remains similar even when take into account their marital status.
4. Purchase values are similar for males and females basis Stay_in_current_city_years. Although, Males buy slightly high value products.

Correlation between categorical variables:#

```
In [40]: numeric_columns = df.select_dtypes(include=['float64', 'int64']).columns

# Compute correlation matrix for numeric columns
correlation_matrix = df[numeric_columns].corr()

# Plot heatmap for correlation matrix
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap="Blues", linewidth=.5)
plt.title('Correlation Matrix Heatmap')
plt.show()
```



## Observations:

1. From the above correlation plot, we can see the correlation is not significant between any pair of variables.

In [41]:
```python
avgamt_gender = df.groupby(['User_ID', 'Gender'])[['Purchase']].sum()
avgamt_gender = avgamt_gender.reset_index()
avgamt_gender
```

Out[41]:

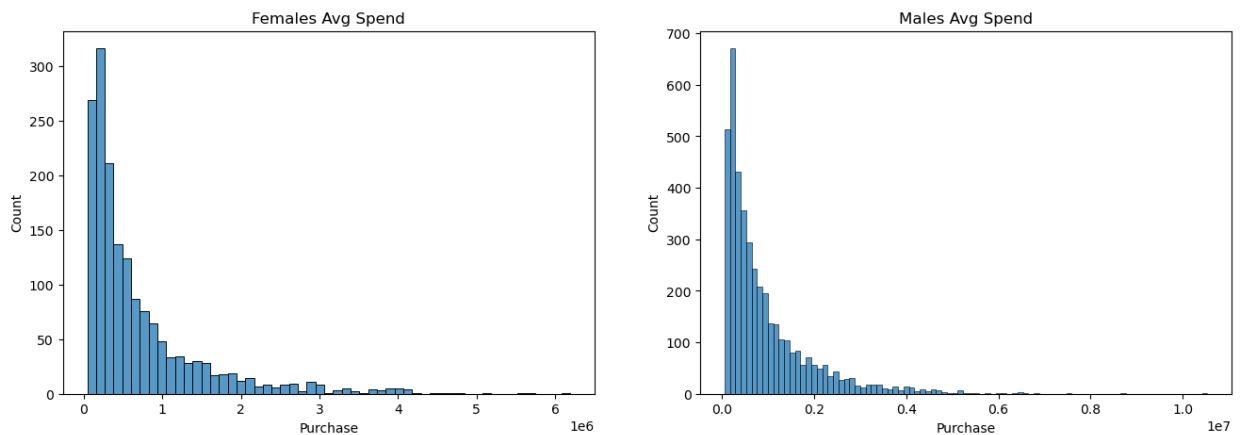|  | User_ID | Gender | Purchase |
|---|---|---|---|
| 0 | 1000001 | F | 334093 |
| 1 | 1000002 | M | 810472 |
| 2 | 1000003 | M | 341635 |
| 3 | 1000004 | M | 206468 |
| 4 | 1000005 | M | 821001 |
| ... | ... | ... | ... |
| 5886 | 1006036 | F | 4116058 |
| 5887 | 1006037 | F | 1119538 |
| 5888 | 1006038 | F | 90034 |
| 5889 | 1006039 | F | 590319 |
| 5890 | 1006040 | M | 1653299 |

5891 rows × 3 columns

In [42]:
```python
# Gender wise count in the entire data
avgamt_gender['Gender'].value_counts()
```

Out[42]:
```
Gender
M    4225
F    1666
Name: count, dtype: int64
```

In [43]:
```python
fig, axs = plt.subplots(nrows=1, ncols=2, figsize=(16,5))

sns.histplot(data=avgamt_gender[avgamt_gender['Gender']=='F']['Purchase'], ax=axs[0]).set_title("Female
sns.histplot(data=avgamt_gender[avgamt_gender['Gender']=='M']['Purchase'], ax=axs[1]).set_title("Males
```

Out[43]: Text(0.5, 1.0, 'Males Avg Spend')



# Observations:

1. Average amount spend by males are higher than females.

In [44]:
```python
avgamt_gender.groupby(['Gender'])[['Purchase']].mean()
```

Out[44]:

|  | Purchase |
|---|---|
| **Gender** | |
| F | 712024.394958 |
| M | 925344.402367 |

```
In [45]: avgamt_gender.groupby(['Gender'])['Purchase'].sum()
```

```
Out[45]: Gender
         F    1186232642
         M    3909580100
         Name: Purchase, dtype: int64
```

## Observations:

1. Average amount for the males is 925344 for the entire population whereas it's much lesser for females(712024).
2. Total amount spend by males is around 4 billion whereas for females it's 1.2 billion.

```
In [46]: avgamt_male = avgamt_gender[avgamt_gender['Gender']=='M']
         avgamt_female = avgamt_gender[avgamt_gender['Gender']=='F']
```

```
In [47]: #Finding the sample(sample size=1000) for avg purchase amount for males and females
         genders = ["M", "F"]

         sample_size = 1000

         num_repitions = 1000
         male_means = []
         female_means = []

         for i in range(num_repitions):
             male_mean = avgamt_male.sample(sample_size, replace=True)['Purchase'].mean()
             female_mean = avgamt_female.sample(sample_size, replace=True)['Purchase'].mean()

             male_means.append(male_mean)
             female_means.append(female_mean)
```
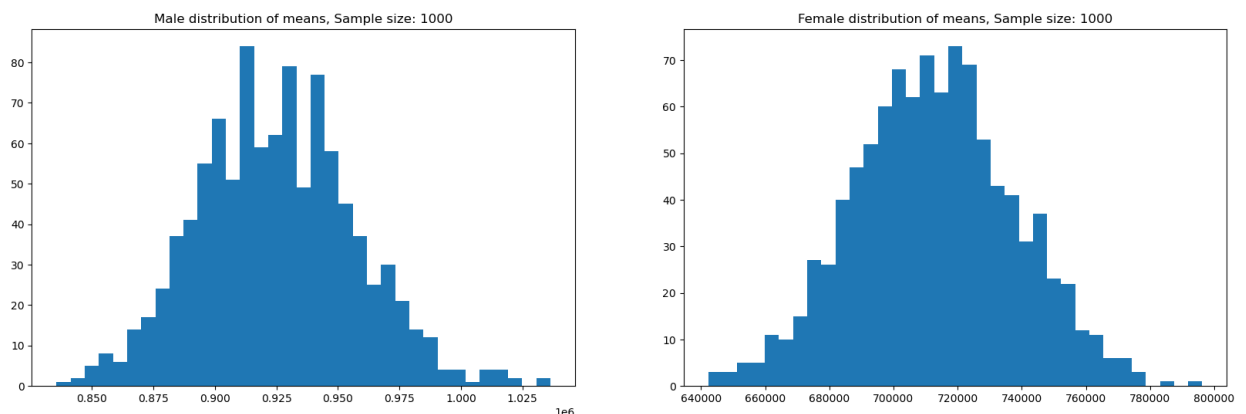
```
In [48]: fig, axis = plt.subplots(nrows=1, ncols=2, figsize=(20, 6))

         axis[0].hist(male_means, bins=35)
         axis[1].hist(female_means, bins=35)
         axis[0].set_title("Male distribution of means, Sample size: 1000")
         axis[1].set_title("Female distribution of means, Sample size: 1000")

         plt.show()
```



## Observations:

1. The means sample seems to be normally distributed for both males and females. Also, we can see the mean of the sample means are closer to the population mean as per central limit theorem.

## Calculating 90% confidence interval for sample size 1000:

In [49]:
```python
#Taking the values for z at 90%, 95% and 99% confidence interval as:
z90=1.645 #90% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

print("Population avg spend amount for Male: {:.2f}".format(avgamt_male['Purchase'].mean()))
print("Population avg spend amount for Female: {:.2f}\n".format(avgamt_female['Purchase'].mean()))

print("Sample avg spend amount for Male: {:.2f}".format(np.mean(male_means)))
print("Sample avg spend amount for Female: {:.2f}\n".format(np.mean(female_means)))

print("Sample std for Male: {:.2f}".format(pd.Series(male_means).std()))
print("Sample std for Female: {:.2f}\n".format(pd.Series(female_means).std()))

print("Sample std error for Male: {:.2f}".format(pd.Series(male_means).std()/np.sqrt(1000)))
print("Sample std error for Female: {:.2f}\n".format(pd.Series(female_means).std()/np.sqrt(1000)))

sample_mean_male=np.mean(male_means)
sample_mean_female=np.mean(female_means)

sample_std_male=pd.Series(male_means).std()
sample_std_female=pd.Series(female_means).std()

sample_std_error_male=sample_std_male/np.sqrt(1000)
sample_std_error_female=sample_std_female/np.sqrt(1000)

Upper_Limit_male=z90*sample_std_error_male + sample_mean_male
Lower_Limit_male=sample_mean_male - z90*sample_std_error_male

Upper_Limit_female=z90*sample_std_error_female + sample_mean_female
Lower_Limit_female=sample_mean_female - z90*sample_std_error_female

print("Male_CI: ",[Lower_Limit_male,Upper_Limit_male])
print("Female_CI: ",[Lower_Limit_female,Upper_Limit_female])
```

```
Population avg spend amount for Male: 925344.40
Population avg spend amount for Female: 712024.39

Sample avg spend amount for Male: 925147.25
Sample avg spend amount for Female: 712817.99

Sample std for Male: 31959.42
Sample std for Female: 25019.06

Sample std error for Male: 1010.65
Sample std error for Female: 791.17

Male_CI:  [923484.7382237061, 926809.7620962937]
Female_CI:  [711516.5117451514, 714119.4680668483]
```

# Observation:

Now using the Confidence interval at 90%, we can say that:

Average amount spend by male customers lie in the range 9,22,940.71 - 9,26,225.18

Average amount spend by female customers lie in range 7,10,425.64 - 7,13,064.55

In [50]:
```python
#Calculating 95% confidence interval for sample size 1000:
```

In [51]:
```python
#Taking the values for z at 90%, 95% and 99% confidence interval as:
z90=1.645 #90% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

print("Population avg spend amount for Male: {:.2f}".format(avgamt_male['Purchase'].mean()))
print("Population avg spend amount for Female: {:.2f}\n".format(avgamt_female['Purchase'].mean()))

print("Sample avg spend amount for Male: {:.2f}".format(np.mean(male_means)))
print("Sample avg spend amount for Female: {:.2f}\n".format(np.mean(female_means)))

print("Sample std for Male: {:.2f}".format(pd.Series(male_means).std()))
print("Sample std for Female: {:.2f}\n".format(pd.Series(female_means).std()))

print("Sample std error for Male: {:.2f}".format(pd.Series(male_means).std()/np.sqrt(1000)))
print("Sample std error for Female: {:.2f}\n".format(pd.Series(female_means).std()/np.sqrt(1000)))

sample_mean_male=np.mean(male_means)
sample_mean_female=np.mean(female_means)

sample_std_male=pd.Series(male_means).std()
sample_std_female=pd.Series(female_means).std()

sample_std_error_male=sample_std_male/np.sqrt(1000)
sample_std_error_female=sample_std_female/np.sqrt(1000)

Upper_Limit_male=z95*sample_std_error_male + sample_mean_male
Lower_Limit_male=sample_mean_male - z95*sample_std_error_male

Upper_Limit_female=z95*sample_std_error_female + sample_mean_female
Lower_Limit_female=sample_mean_female - z95*sample_std_error_female

print("Male_CI: ",[Lower_Limit_male,Upper_Limit_male])
print("Female_CI: ",[Lower_Limit_female,Upper_Limit_female])
```

```
Population avg spend amount for Male: 925344.40
Population avg spend amount for Female: 712024.39

Sample avg spend amount for Male: 925147.25
Sample avg spend amount for Female: 712817.99

Sample std for Male: 31959.42
Sample std for Female: 25019.06

Sample std error for Male: 1010.65
Sample std error for Female: 791.17

Male_CI:  [923166.384874203, 927128.1154457969]
Female_CI:  [711267.2925228612, 714368.6872891386]
```

# Observation:

Now using the Confidence interval at 95%, we can say that:

Average amount spend by male customers lie in the range 9,22,626.24 - 9,26,539.65

Average amount spend by female customers lie in range 7,10,172.98 - 7,13,317.21

In [52]:
```python
#Calculating 99% confidence interval for sample size 1000:
```

```
In [53]: #Taking the values for z at 90%, 95% and 99% confidence interval as:
         z90=1.645 #90% Confidence Interval
         z95=1.960 #95% Confidence Interval
         z99=2.576 #99% Confidence Interval

         print("Population avg spend amount for Male: {:.2f}".format(avgamt_male['Purchase'].mean()))
         print("Population avg spend amount for Female: {:.2f}\n".format(avgamt_female['Purchase'].mean()))

         print("Sample avg spend amount for Male: {:.2f}".format(np.mean(male_means)))
         print("Sample avg spend amount for Female: {:.2f}\n".format(np.mean(female_means)))

         print("Sample std for Male: {:.2f}".format(pd.Series(male_means).std()))
         print("Sample std for Female: {:.2f}\n".format(pd.Series(female_means).std()))

         print("Sample std error for Male: {:.2f}".format(pd.Series(male_means).std()/np.sqrt(1000)))
         print("Sample std error for Female: {:.2f}\n".format(pd.Series(female_means).std()/np.sqrt(1000)))

         sample_mean_male=np.mean(male_means)
         sample_mean_female=np.mean(female_means)

         sample_std_male=pd.Series(male_means).std()
         sample_std_female=pd.Series(female_means).std()

         sample_std_error_male=sample_std_male/np.sqrt(1000)
         sample_std_error_female=sample_std_female/np.sqrt(1000)

         Upper_Limit_male=z99*sample_std_error_male + sample_mean_male
         Lower_Limit_male=sample_mean_male - z99*sample_std_error_male

         Upper_Limit_female=z99*sample_std_error_female + sample_mean_female
         Lower_Limit_female=sample_mean_female - z99*sample_std_error_female

         print("Male_CI: ",[Lower_Limit_male,Upper_Limit_male])
         print("Female_CI: ",[Lower_Limit_female,Upper_Limit_female])
```

```
Population avg spend amount for Male: 925344.40
Population avg spend amount for Female: 712024.39

Sample avg spend amount for Male: 925147.25
Sample avg spend amount for Female: 712817.99

Sample std for Male: 31959.42
Sample std for Female: 25019.06

Sample std error for Male: 1010.65
Sample std error for Female: 791.17

Male_CI:  [922543.8272129525, 927750.6731070473]
Female_CI:  [710779.9304881605, 714856.0493238393]
```

# Observation:

Now using the Confidence interval at 99%, we can say that:

Average amount spend by male customers lie in the range 9,22,011.28 - 9,27,154.61

Average amount spend by female customers lie in range 7,09,678.88 - 7,13,811.31

```
In [54]: #Calculating 90% confidence interval for sample size 1500:
```

In [55]:
```python
#Finding the sample(sample size=1000) avg purchase amount for males and females
genders = ["M", "F"]

sample_size = 1500

num_repitions = 1000
male_means = []
female_means = []

for i in range(num_repitions):
    male_mean = avgamt_male.sample(sample_size, replace=True)['Purchase'].mean()
    female_mean = avgamt_female.sample(sample_size, replace=True)['Purchase'].mean()

    male_means.append(male_mean)
    female_means.append(female_mean)

#Taking the values for z at 90%, 95% and 99% confidence interval as:
z90=1.645 #90% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

print("Population avg spend amount for Male: {:.2f}".format(avgamt_male['Purchase'].mean()))
print("Population avg spend amount for Female: {:.2f}\n".format(avgamt_female['Purchase'].mean()))

print("Sample avg spend amount for Male: {:.2f}".format(np.mean(male_means)))
print("Sample avg spend amount for Female: {:.2f}\n".format(np.mean(female_means)))

print("Sample std for Male: {:.2f}".format(pd.Series(male_means).std()))
print("Sample std for Female: {:.2f}\n".format(pd.Series(female_means).std()))

print("Sample std error for Male: {:.2f}".format(pd.Series(male_means).std()/np.sqrt(1500)))
print("Sample std error for Female: {:.2f}\n".format(pd.Series(female_means).std()/np.sqrt(1500)))

sample_mean_male=np.mean(male_means)
sample_mean_female=np.mean(female_means)

sample_std_male=pd.Series(male_means).std()
sample_std_female=pd.Series(female_means).std()

sample_std_error_male=sample_std_male/np.sqrt(1500)
sample_std_error_female=sample_std_female/np.sqrt(1500)

Upper_Limit_male=z90*sample_std_error_male + sample_mean_male
Lower_Limit_male=sample_mean_male - z90*sample_std_error_male

Upper_Limit_female=z90*sample_std_error_female + sample_mean_female
Lower_Limit_female=sample_mean_female - z90*sample_std_error_female

print("Male_CI: ",[Lower_Limit_male,Upper_Limit_male])
print("Female_CI: ",[Lower_Limit_female,Upper_Limit_female])
```

```
Population avg spend amount for Male: 925344.40
Population avg spend amount for Female: 712024.39

Sample avg spend amount for Male: 924731.80
Sample avg spend amount for Female: 713038.16

Sample std for Male: 25795.00
Sample std for Female: 21039.44

Sample std error for Male: 666.02
Sample std error for Female: 543.24

Male_CI:  [923636.1939735838, 925827.4134317496]
Female_CI:  [712144.5326646747, 713931.7787353253]
```

# Observation:

Now using the Confidence interval at 90%, we can say that:

Average amount spend by male customers lie in the range 9,24,177.41 - 9,26,318.90

Average amount spend by female customers lie in range 7,11,187.27 - 7,12,971.67

By increasing the sample size we can see confidence interval is more closer to the population mean.

In [56]: ```python
#Calculating 95% confidence interval for sample size 1500:
```

In [57]: ```python
print("Population avg spend amount for Male: {:.2f}".format(avgamt_male['Purchase'].mean()))
print("Population avg spend amount for Female: {:.2f}\n".format(avgamt_female['Purchase'].mean()))

print("Sample avg spend amount for Male: {:.2f}".format(np.mean(male_means)))
print("Sample avg spend amount for Female: {:.2f}\n".format(np.mean(female_means)))

print("Sample std for Male: {:.2f}".format(pd.Series(male_means).std()))
print("Sample std for Female: {:.2f}\n".format(pd.Series(female_means).std()))

print("Sample std error for Male: {:.2f}".format(pd.Series(male_means).std()/np.sqrt(1500)))
print("Sample std error for Female: {:.2f}\n".format(pd.Series(female_means).std()/np.sqrt(1500)))

sample_mean_male=np.mean(male_means)
sample_mean_female=np.mean(female_means)

sample_std_male=pd.Series(male_means).std()
sample_std_female=pd.Series(female_means).std()

sample_std_error_male=sample_std_male/np.sqrt(1500)
sample_std_error_female=sample_std_female/np.sqrt(1500)

Upper_Limit_male=z95*sample_std_error_male + sample_mean_male
Lower_Limit_male=sample_mean_male - z95*sample_std_error_male

Upper_Limit_female=z95*sample_std_error_female + sample_mean_female
Lower_Limit_female=sample_mean_female - z95*sample_std_error_female

print("Male_CI: ",[Lower_Limit_male,Upper_Limit_male])
print("Female_CI: ",[Lower_Limit_female,Upper_Limit_female])
```

```
Population avg spend amount for Male: 925344.40
Population avg spend amount for Female: 712024.39

Sample avg spend amount for Male: 924731.80
Sample avg spend amount for Female: 713038.16

Sample std for Male: 25795.00
Sample std for Female: 21039.44

Sample std error for Male: 666.02
Sample std error for Female: 543.24

Male_CI:  [923426.396365887, 926037.2110394463]
Female_CI:  [711973.413360038, 714102.898039962]
```

# Observation:

Now using the Confidence interval at 95%, we can say that:

Average amount spend by male customers lie in the range 9,23,972.41 - 9,26,523.93

Average amount spend by female customers lie in range 7,11,016.42 - 7,13,142.51

By increasing the sample size we can see confidence interval is more closer to the population mean.

In [58]: ```python
#Calculating 99% confidence interval for sample size 1500:
```

```
In [59]: print("Population avg spend amount for Male: {:.2f}".format(avgamt_male['Purchase'].mean()))
         print("Population avg spend amount for Female: {:.2f}\n".format(avgamt_female['Purchase'].mean()))

         print("Sample avg spend amount for Male: {:.2f}".format(np.mean(male_means)))
         print("Sample avg spend amount for Female: {:.2f}\n".format(np.mean(female_means)))

         print("Sample std for Male: {:.2f}".format(pd.Series(male_means).std()))
         print("Sample std for Female: {:.2f}\n".format(pd.Series(female_means).std()))

         print("Sample std error for Male: {:.2f}".format(pd.Series(male_means).std()/np.sqrt(1500)))
         print("Sample std error for Female: {:.2f}\n".format(pd.Series(female_means).std()/np.sqrt(1500)))

         sample_mean_male=np.mean(male_means)
         sample_mean_female=np.mean(female_means)

         sample_std_male=pd.Series(male_means).std()
         sample_std_female=pd.Series(female_means).std()

         sample_std_error_male=sample_std_male/np.sqrt(1500)
         sample_std_error_female=sample_std_female/np.sqrt(1500)

         Upper_Limit_male=z99*sample_std_error_male + sample_mean_male
         Lower_Limit_male=sample_mean_male - z99*sample_std_error_male

         Upper_Limit_female=z99*sample_std_error_female + sample_mean_female
         Lower_Limit_female=sample_mean_female - z99*sample_std_error_female

         print("Male_CI: ",[Lower_Limit_male,Upper_Limit_male])
         print("Female_CI: ",[Lower_Limit_female,Upper_Limit_female])
```

```
Population avg spend amount for Male: 925344.40
Population avg spend amount for Female: 712024.39

Sample avg spend amount for Male: 924731.80
Sample avg spend amount for Female: 713038.16

Sample std for Male: 25795.00
Sample std for Female: 21039.44

Sample std error for Male: 666.02
Sample std error for Female: 543.24

Male_CI:  [923016.1254886135, 926447.4819167198]
Female_CI:  [711638.7800531927, 714437.5313468073]
```

## Observation:

Now using the Confidence interval at 99%, we can say that:

Average amount spend by male customers lie in the range 923571.42 - 926924.89

Average amount spend by female customers lie in range 710682.32 - 713476.61

By increasing the sample size we can see confidence interval is more closer to the population mean.

```
In [60]: #CLT and Confidence interval considering marital status:
```

In [61]:
```python
avg_Marital = df.groupby(['User_ID', 'Marital_Status'])[['Purchase']].sum()
avg_Marital = avg_Marital.reset_index()

avgamt_married = avg_Marital[avg_Marital['Marital_Status']==1]
avgamt_single = avg_Marital[avg_Marital['Marital_Status']==0]

sample_size = 1000
num_repitions = 1000
married_means = []
single_means = []

for i in range(num_repitions):
    avg_married = avg_Marital[avg_Marital['Marital_Status']==1].sample(sample_size, replace=True)['Purcl
    avg_single = avg_Marital[avg_Marital['Marital_Status']==0].sample(sample_size, replace=True)['Purcha
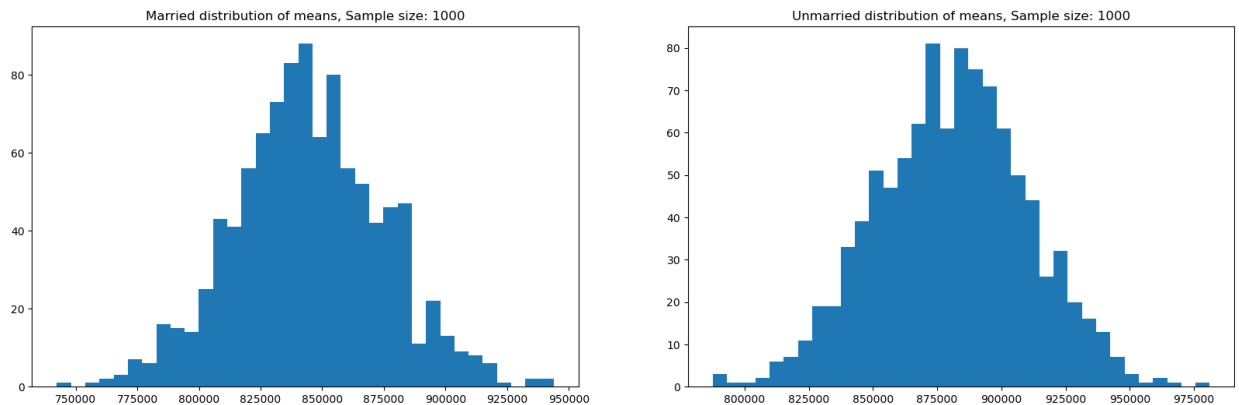
    married_means.append(avg_married)
    single_means.append(avg_single)


fig, axis = plt.subplots(nrows=1, ncols=2, figsize=(20, 6))

axis[0].hist(married_means, bins=35)
axis[1].hist(single_means, bins=35)
axis[0].set_title("Married distribution of means, Sample size: 1000")
axis[1].set_title("Unmarried distribution of means, Sample size: 1000")

plt.show()
```



## Observations:

1. The means sample seems to be normally distributed for both married and singles. Also, we can see the mean of the sample means are closer to the population mean as per central limit theorem.

In [62]:
```python
avg_Marital['Marital_Status'].value_counts()
```

Out[62]:
```
Marital_Status
0    3417
1    2474
Name: count, dtype: int64
```

In [63]:
```python
#Calculating 90% confidence interval for avg expenses for married/single for sample size 1000:
```

```
In [64]: #Taking the values for z at 90%, 95% and 99% confidence interval as:
         z90=1.645 #90% Confidence Interval
         z95=1.960 #95% Confidence Interval
         z99=2.576 #99% Confidence Interval

         print("Population avg spend amount for Married: {:.2f}".format(avgamt_married['Purchase'].mean()))
         print("Population avg spend amount for Single: {:.2f}\n".format(avgamt_single['Purchase'].mean()))

         print("Sample avg spend amount for Married: {:.2f}".format(np.mean(married_means)))
         print("Sample avg spend amount for Single: {:.2f}\n".format(np.mean(single_means)))

         print("Sample std for Married: {:.2f}".format(pd.Series(married_means).std()))
         print("Sample std for Single: {:.2f}\n".format(pd.Series(single_means).std()))

         print("Sample std error for Married: {:.2f}".format(pd.Series(married_means).std()/np.sqrt(1000)))
         print("Sample std error for Single: {:.2f}\n".format(pd.Series(single_means).std()/np.sqrt(1000)))

         sample_mean_married=np.mean(married_means)
         sample_mean_single=np.mean(single_means)

         sample_std_married=pd.Series(married_means).std()
         sample_std_single=pd.Series(single_means).std()

         sample_std_error_married=sample_std_married/np.sqrt(1000)
         sample_std_error_single=sample_std_single/np.sqrt(1000)

         Upper_Limit_married=z90*sample_std_error_male + sample_mean_married
         Lower_Limit_married=sample_mean_married - z90*sample_std_error_married

         Upper_Limit_single=z90*sample_std_error_single + sample_mean_single
         Lower_Limit_single=sample_mean_single - z90*sample_std_error_single

         print("Married_CI: ",[Lower_Limit_married,Upper_Limit_married])
         print("Single_CI: ",[Lower_Limit_single,Upper_Limit_single])
```

```
Population avg spend amount for Married: 843526.80
Population avg spend amount for Single: 880575.78

Sample avg spend amount for Married: 844633.80
Sample avg spend amount for Single: 880840.44

Sample std for Married: 30152.35
Sample std for Single: 29468.78

Sample std error for Married: 953.50
Sample std error for Single: 931.88

Married_CI:  [843065.2941878818, 845729.4131140829]
Single_CI:  [879307.4895850717, 882373.3903809284]
```

```
In [65]: #Calculating 95% confidence interval for avg expenses for married/single for sample size 1000:
```

In [66]:
```python
#Taking the values for z at 90%, 95% and 99% confidence interval as:
z90=1.645 #90% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

print("Population avg spend amount for Married: {:.2f}".format(avgamt_married['Purchase'].mean()))
print("Population avg spend amount for Single: {:.2f}\n".format(avgamt_single['Purchase'].mean()))

print("Sample avg spend amount for Married: {:.2f}".format(np.mean(married_means)))
print("Sample avg spend amount for Single: {:.2f}\n".format(np.mean(single_means)))

print("Sample std for Married: {:.2f}".format(pd.Series(married_means).std()))
print("Sample std for Single: {:.2f}\n".format(pd.Series(single_means).std()))

print("Sample std error for Married: {:.2f}".format(pd.Series(married_means).std()/np.sqrt(1000)))
print("Sample std error for Single: {:.2f}\n".format(pd.Series(single_means).std()/np.sqrt(1000)))

sample_mean_married=np.mean(married_means)
sample_mean_single=np.mean(single_means)

sample_std_married=pd.Series(married_means).std()
sample_std_single=pd.Series(single_means).std()

sample_std_error_married=sample_std_married/np.sqrt(1000)
sample_std_error_single=sample_std_single/np.sqrt(1000)

Upper_Limit_married=z95*sample_std_error_male + sample_mean_married
Lower_Limit_married=sample_mean_married - z95*sample_std_error_married

Upper_Limit_single=z95*sample_std_error_single + sample_mean_single
Lower_Limit_single=sample_mean_single - z95*sample_std_error_single

print("Married_CI: ",[Lower_Limit_married,Upper_Limit_married])
print("Single_CI: ",[Lower_Limit_single,Upper_Limit_single])
```

```
Population avg spend amount for Married: 843526.80
Population avg spend amount for Single: 880575.78

Sample avg spend amount for Married: 844633.80
Sample avg spend amount for Single: 880840.44

Sample std for Married: 30152.35
Sample std for Single: 29468.78

Sample std error for Married: 953.50
Sample std error for Single: 931.88

Married_CI:  [842764.9413629017, 845939.2107217796]
Single_CI:  [879013.9458918514, 882666.9340741487]
```

In [67]:
```python
#Calculating 99% confidence interval for avg expenses for married/single for sample size 1000:
```

```python
In [68]: #Taking the values for z at 90%, 95% and 99% confidence interval as:
         z90=1.645 #90% Confidence Interval
         z95=1.960 #95% Confidence Interval
         z99=2.576 #99% Confidence Interval

         print("Population avg spend amount for Married: {:.2f}".format(avgamt_married['Purchase'].mean()))
         print("Population avg spend amount for Single: {:.2f}\n".format(avgamt_single['Purchase'].mean()))

         print("Sample avg spend amount for Married: {:.2f}".format(np.mean(married_means)))
         print("Sample avg spend amount for Single: {:.2f}\n".format(np.mean(single_means)))

         print("Sample std for Married: {:.2f}".format(pd.Series(married_means).std()))
         print("Sample std for Single: {:.2f}\n".format(pd.Series(single_means).std()))

         print("Sample std error for Married: {:.2f}".format(pd.Series(married_means).std()/np.sqrt(1000)))
         print("Sample std error for Single: {:.2f}\n".format(pd.Series(single_means).std()/np.sqrt(1000)))

         sample_mean_married=np.mean(married_means)
         sample_mean_single=np.mean(single_means)

         sample_std_married=pd.Series(married_means).std()
         sample_std_single=pd.Series(single_means).std()

         sample_std_error_married=sample_std_married/np.sqrt(1000)
         sample_std_error_single=sample_std_single/np.sqrt(1000)

         Upper_Limit_married=z99*sample_std_error_male + sample_mean_married
         Lower_Limit_married=sample_mean_married - z99*sample_std_error_married

         Upper_Limit_single=z99*sample_std_error_single + sample_mean_single
         Lower_Limit_single=sample_mean_single - z99*sample_std_error_single

         print("Married_CI: ",[Lower_Limit_married,Upper_Limit_married])
         print("Single_CI: ",[Lower_Limit_single,Upper_Limit_single])
```

```
Population avg spend amount for Married: 843526.80
Population avg spend amount for Single: 880575.78

Sample avg spend amount for Married: 844633.80
Sample avg spend amount for Single: 880840.44

Sample std for Married: 30152.35
Sample std for Single: 29468.78

Sample std error for Married: 953.50
Sample std error for Single: 931.88

Married_CI:  [842177.5847273851, 846349.4815990531]
Single_CI:  [878439.9048917762, 883240.9750742239]
```

# Observation:

For married and singles, it can be seen with larger sample size the sample mean gets closer to tthe population mean. And at greater confidence interval, the range increases.

```python
In [69]: avgamt_age = df.groupby(['User_ID', 'Age'])[['Purchase']].sum()
         avgamt_age = avgamt_age.reset_index()

         avgamt_age['Age'].value_counts()
```

```
Out[69]: Age
         26-35    2053
         36-45    1167
         18-25    1069
         46-50     531
         51-55     481
         55+       372
         0-17      218
         Name: count, dtype: int64
```

```
In [70]: sample_size = 200
         num_repitions = 1000

         all_sample_means = {}

         age_intervals = ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']
         for i in age_intervals:
             all_sample_means[i] = []

         for i in age_intervals:
             for j in range(num_repitions):

                 mean = avgamt_age[avgamt_age['Age']==i].sample(sample_size, replace=True)['Purchase'].mean()
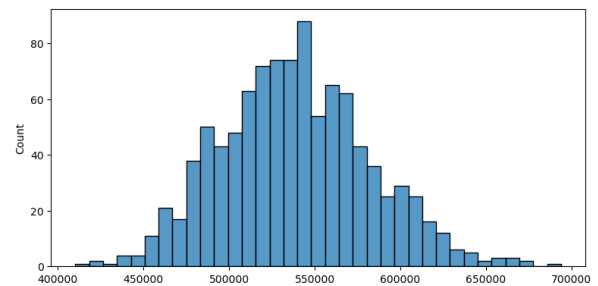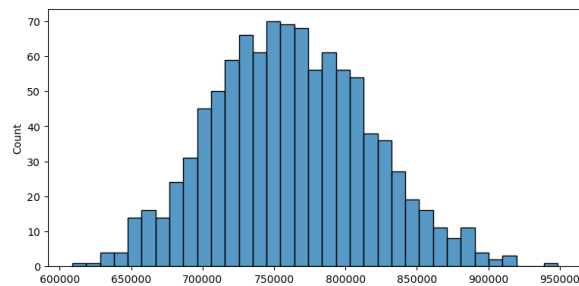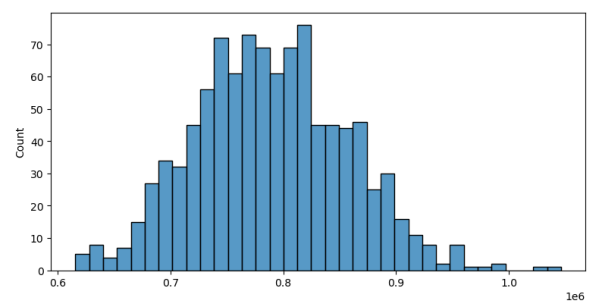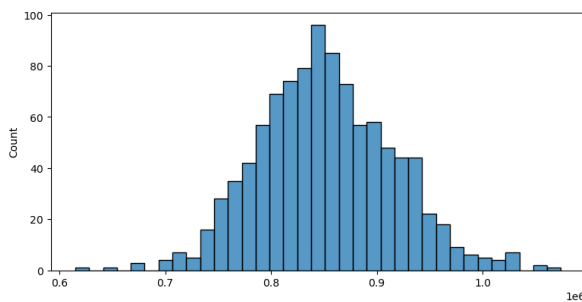                 all_sample_means[i].append(mean)
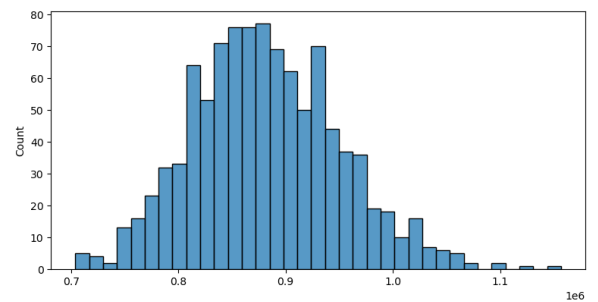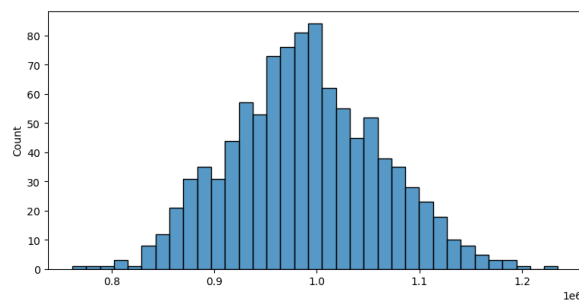


         fig, axis = plt.subplots(nrows=3, ncols=2, figsize=(20, 15))
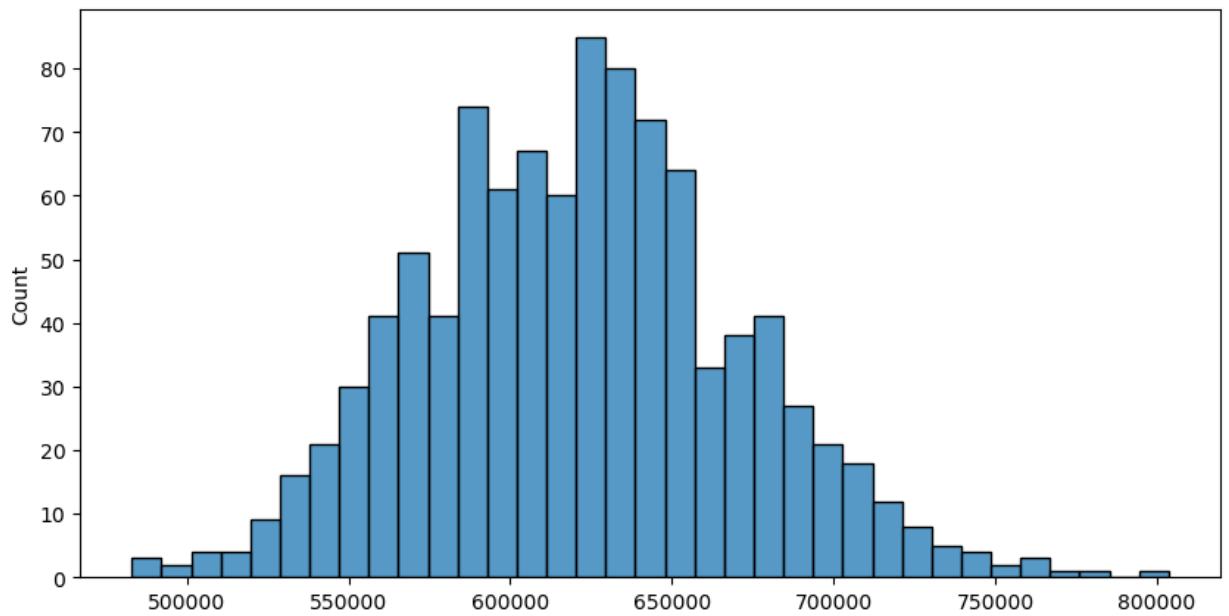
         sns.histplot(all_sample_means['26-35'],bins=35,ax=axis[0,0])
         sns.histplot(all_sample_means['36-45'],bins=35,ax=axis[0,1])
         sns.histplot(all_sample_means['18-25'],bins=35,ax=axis[1,0])
         sns.histplot(all_sample_means['46-50'],bins=35,ax=axis[1,1])
         sns.histplot(all_sample_means['51-55'],bins=35,ax=axis[2,0])
         sns.histplot(all_sample_means['55+'],bins=35,ax=axis[2,1])

         plt.show()

         plt.figure(figsize=(10, 5))
         sns.histplot(all_sample_means['0-17'],bins=35)
         plt.show()
```

## Observations:

1. The means sample seems to be normally distributed for all age groups. Also, we can see the mean of the sample means are closer to the population mean as per central limit theorem.

In [71]: ```#Calculating 90% confidence interval for avg expenses for different age groups for sample size 200:```

In [72]:
```python
z90=1.645 #90% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

sample_size = 200
num_repitions = 1000

all_population_means={}
all_sample_means = {}

age_intervals = ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']
for i in age_intervals:
    all_sample_means[i] = []
    all_population_means[i]=[]
    population_mean=avgamt_age[avgamt_age['Age']==i]['Purchase'].mean()
    all_population_means[i].append(population_mean)

print("All age group population mean: \n", all_population_means)
print("\n")

for i in age_intervals:
    for j in range(num_repitions):

        mean = avgamt_age[avgamt_age['Age']==i].sample(sample_size, replace=True)['Purchase'].mean()
        all_sample_means[i].append(mean)


for val in ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']:

    new_df = avgamt_age[avgamt_age['Age']==val]

    std_error = z90*new_df['Purchase'].std()/np.sqrt(len(new_df))
    sample_mean = new_df['Purchase'].mean()
    lower_lim = sample_mean - std_error
    upper_lim = sample_mean + std_error

    print("For age {} confidence interval of means: ({:.2f}, {:.2f})".format(val, lower_lim, upper_lim)
```

```
All age group population mean:
 {'26-35': [989659.3170969313], '36-45': [879665.7103684661], '18-25': [854863.119738073], '46-50': [7
92548.7815442561], '51-55': [763200.9230769231], '55+': [539697.2446236559], '0-17': [618867.811926605
5]}


For age 26-35 confidence interval of means: (952206.28, 1027112.35)
For age 36-45 confidence interval of means: (832398.89, 926932.53)
For age 18-25 confidence interval of means: (810187.65, 899538.59)
For age 46-50 confidence interval of means: (726209.00, 858888.57)
For age 51-55 confidence interval of means: (703772.36, 822629.48)
For age 55+ confidence interval of means: (487032.92, 592361.57)
For age 0-17 confidence interval of means: (542320.46, 695415.16)
```

In [73]:
```python
#Calculating 95% confidence interval for avg expenses for different age groups for sample size 200:
```

In [74]:
```python
z90=1.645 #90% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

sample_size = 200
num_repitions = 1000

all_means = {}

age_intervals = ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']
for i in age_intervals:
    all_means[i] = []

for i in age_intervals:
    for j in range(num_repitions):
        mean = avgamt_age[avgamt_age['Age']==i].sample(sample_size, replace=True)['Purchase'].mean()
        all_means[i].append(mean)
for val in ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']:

    new_df = avgamt_age[avgamt_age['Age']==val]

    std_error = z95*new_df['Purchase'].std()/np.sqrt(len(new_df))
    sample_mean = new_df['Purchase'].mean()
    lower_lim = sample_mean - std_error
    upper_lim = sample_mean + std_error

    print("For age {} confidence interval of means: ({:.2f}, {:.2f})".format(val, lower_lim, upper_lim)
```

```
For age 26-35 confidence interval of means: (945034.42, 1034284.21)
For age 36-45 confidence interval of means: (823347.80, 935983.62)
For age 18-25 confidence interval of means: (801632.78, 908093.46)
For age 46-50 confidence interval of means: (713505.63, 871591.93)
For age 51-55 confidence interval of means: (692392.43, 834009.42)
For age 55+ confidence interval of means: (476948.26, 602446.23)
For age 0-17 confidence interval of means: (527662.46, 710073.17)
```

## Calculating 99% confidence interval for avg expenses for different age groups for sample size 200:

In [75]:
```python
z90=1.645 #90% Confidence Interval
z95=1.960 #95% Confidence Interval
z99=2.576 #99% Confidence Interval

sample_size = 200
num_repitions = 1000

all_means = {}

age_intervals = ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']
for i in age_intervals:
    all_means[i] = []

for i in age_intervals:
    for j in range(num_repitions):
        mean = avgamt_age[avgamt_age['Age']==i].sample(sample_size, replace=True)['Purchase'].mean()
        all_means[i].append(mean)
for val in ['26-35', '36-45', '18-25', '46-50', '51-55', '55+', '0-17']:

    new_df = avgamt_age[avgamt_age['Age']==val]

    std_error = z99*new_df['Purchase'].std()/np.sqrt(len(new_df))
    sample_mean = new_df['Purchase'].mean()
    lower_lim = sample_mean - std_error
    upper_lim = sample_mean + std_error

    print("For age {} confidence interval of means: ({:.2f}, {:.2f})".format(val, lower_lim, upper_lim)
```

```
For age 26-35 confidence interval of means: (931009.46, 1048309.18)
For age 36-45 confidence interval of means: (805647.89, 953683.53)
For age 18-25 confidence interval of means: (784903.24, 924823.00)
For age 46-50 confidence interval of means: (688663.50, 896434.06)
For age 51-55 confidence interval of means: (670138.33, 856263.52)
For age 55+ confidence interval of means: (457227.15, 622167.34)
For age 0-17 confidence interval of means: (498997.92, 738737.71)
```

# Observation:

1. We can see the sample means are closer to the population mean for the differnt age groups. And, with greater confidence interval we have the upper limit and lower limit range increases. As we have seen for gender and marital status, by increasing the sample size we can have the mean of the sample means closer to the population.

Observations:
1. Male customers are high in number compared to Female .
2. Purchase distribution is similar across Genders.
3. Purchase Distribution is simlar across Genders irrespective of city category.
4. Ratio of Male to Female participatoin is very much similar across all city categories.
5. Distribution of Age is similar in both Genders.
6. [4,0,7,1,17] are top 5 occupations of purchasers.

Recommendations:
1. As the participation of males is high, Male customers can be targeted for better Purchases.
2. Of all Males Customers, those of age group of 26-36 are promising customers, these audiences have high probability of purchasing.
3. Insted of channeling resources and effort over all occupational groups. Resource allocaton should be based on their frequency of purchase.
4. Resource allocation should be done more on bringing the low frequency groups on to platform. which in turn increase the overall purchase.
5. Better Category Specific incentives should be provided based on their category to attract more Users.

In [ ]: