In [4]:
```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

import os
```

In [5]:
```python
df = pd.read_csv(r"C:\Users\Keremane\OneDrive\Desktop\Scaler Docs\Data Sources\Python Pandas\walmart_data.csv")
df.head()
```

Out[5]:

| | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Years | Marital_Status | Product_Category | Purchase |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1000001 | P00069042 | F | 0-17 | 10 | A | 2 | 0 | 3 | 8370 |
| 1 | 1000001 | P00248942 | F | 0-17 | 10 | A | 2 | 0 | 1 | 15200 |
| 2 | 1000001 | P00087842 | F | 0-17 | 10 | A | 2 | 0 | 12 | 1422 |
| 3 | 1000001 | P00085442 | F | 0-17 | 10 | A | 2 | 0 | 12 | 1057 |
| 4 | 1000002 | P00285442 | M | 55+ | 16 | C | 4+ | 0 | 8 | 7969 |

In [6]:
```python
print(f"Number of rows: {df.shape[0]:,} \nNumber of columns: {df.shape[1]}")
```

```
Number of rows: 550,068
Number of columns: 10
```

In [7]:
```python
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
```

In [9]:
```python
cols = ['Occupation', 'Marital_Status', 'Product_Category']
df[cols] = df[cols].astype('object')
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  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(2), object(8)
memory usage: 42.0+ MB
```

In [10]:
```python
df.dtypes
```

Out[10]:
```
User_ID                        int64
Product_ID                    object
Gender                        object
Age                           object
Occupation                    object
City_Category                 object
Stay_In_Current_City_Years    object
Marital_Status                object
Product_Category              object
Purchase                       int64
dtype: object
```

In [11]: `df.describe()`

Out[11]:

|  | User_ID | Purchase |
|---|---|---|
| count | 5.500680e+05 | 550068.000000 |
| mean | 1.003029e+06 | 9263.968713 |
| std | 1.727592e+03 | 5023.065394 |
| min | 1.000001e+06 | 12.000000 |
| 25% | 1.001516e+06 | 5823.000000 |
| 50% | 1.003077e+06 | 8047.000000 |
| 75% | 1.004478e+06 | 12054.000000 |
| max | 1.006040e+06 | 23961.000000 |

In [12]:
```
# checking null values
df.isnull().sum()
```

Out[12]:
```
User_ID                       0
Product_ID                    0
Gender                        0
Age                           0
Occupation                    0
City_Category                 0
Stay_In_Current_City_Years    0
Marital_Status                0
Product_Category              0
Purchase                      0
dtype: int64
```

### Observations

There are no missing values in the dataset.

In [13]: `df['User_ID'].nunique()`

Out[13]: 5891

In [14]: `df['Product_ID'].nunique()`

Out[14]: 3631

```python
categorical_cols = ['Gender', 'Age', 'Occupation', 'City_Category', 'Stay_In_Current_City_Years', 'Marital_Status', 'Product_
df[categorical_cols].melt().groupby(['variable', 'value'])[['value']].count()/len(df)
```

Out[15]:

| variable | value | 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 |
| Occupation | 0 | 0.126599 |
| | 1 | 0.086218 |
| | 2 | 0.048336 |
| | 3 | 0.032087 |
| | 4 | 0.131453 |
| | 5 | 0.022137 |
| | 6 | 0.037005 |
| | 7 | 0.107501 |
| | 8 | 0.002811 |
| | 9 | 0.011437 |
| | 10 | 0.023506 |
| | 11 | 0.021063 |
| | 12 | 0.056682 |
| | 13 | 0.014049 |
| | 14 | 0.049647 |
| | 15 | 0.022115 |
| | 16 | 0.046123 |
| | 17 | 0.072796 |
| | 18 | 0.012039 |
| | 19 | 0.015382 |
| | 20 | 0.061014 |
| Product_Category | 1 | 0.255201 |
| | 2 | 0.043384 |
| | 3 | 0.036746 |
| | 4 | 0.021366 |
| | 5 | 0.274390 |
| | 6 | 0.037206 |
| | 7 | 0.006765 |
| | 8 | 0.207111 |
| | 9 | 0.000745 |
| | 10 | 0.009317 |
| | 11 | 0.044153 |
| | 12 | 0.007175 |
| | 13 | 0.010088 |
| | 14 | 0.002769 |
| | 15 | 0.011435 |
| | 16 | 0.017867 |
| | 17 | 0.001051 |
| | 18 | 0.005681 |
| | 19 | 0.002914 |
| | 20 | 0.004636 |

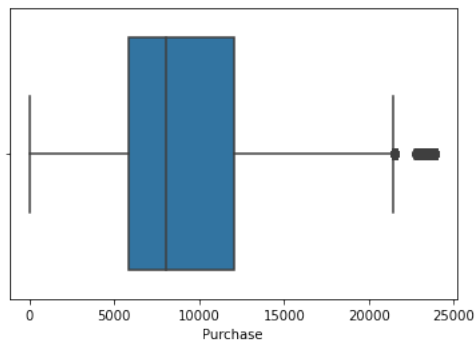|  |  | value |
| --- | --- | --- |
| **variable** | **value** |  |
| **Stay_In_Current_City_Years** | **0** | 0.135252 |
|  | **1** | 0.352358 |
|  | **2** | 0.185137 |
|  | **3** | 0.173224 |
|  | **4+** | 0.154028 |

## Observations

- 80% of the users are between the age 18-50 (40%: 26-35, 18%: 18-25, 20%: 36-45)
- 75% of the users are Male and 25% are Female
- 60% Single, 40% Married
- 35% Staying in the city from 1 year, 18% from 2 years, 17% from 3 years
- Total of 20 product categories are there
- There are 20 differnent types of occupations in the city

# Univariate Analysis

In [17]:
```python
sns.histplot(data=df, x='Purchase', kde=True)
plt.show()
```



In [18]:
```python
sns.boxplot(data=df, x='Purchase', orient='h')
plt.show()
```
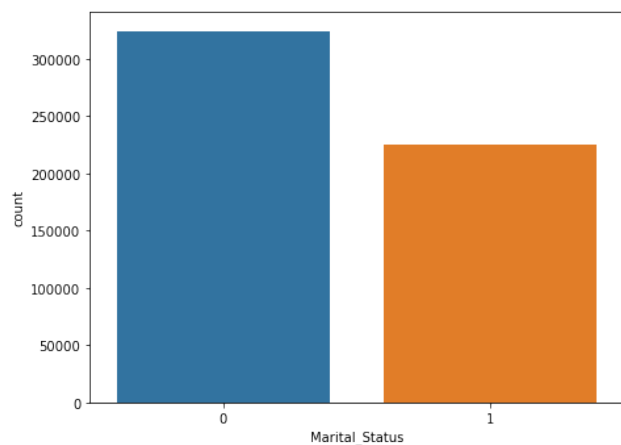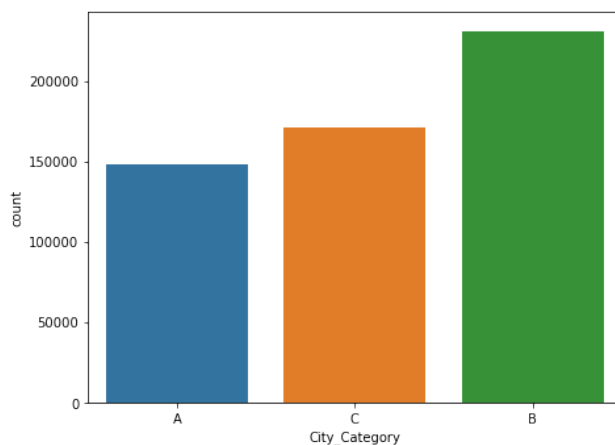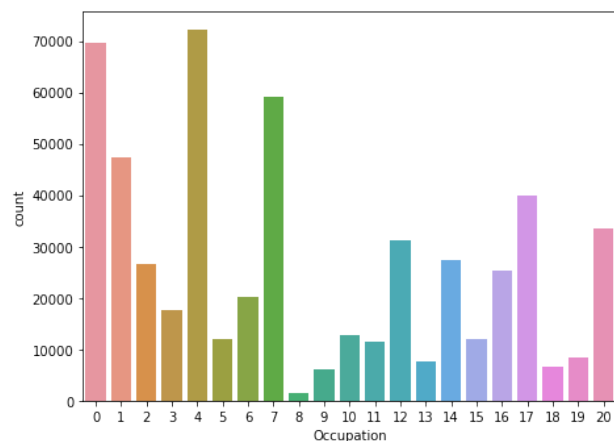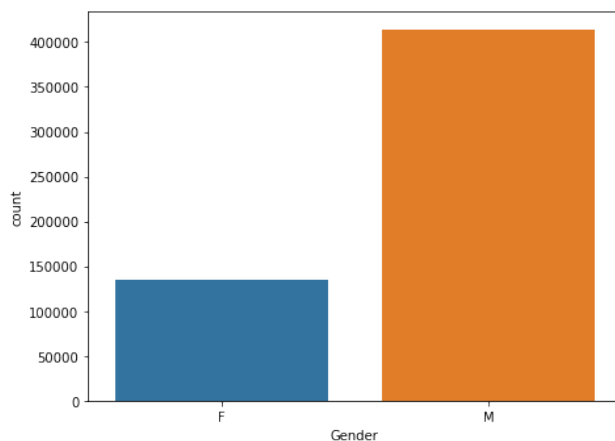


## Observation

Purchase is having outliers

```
In [22]:  categorical_cols = ['Gender', 'Occupation','City_Category','Marital_Status','Product_Category']

          fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(16, 12))
          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()

          plt.figure(figsize=(16, 8))
          sns.countplot(data=df, x='Product_Category')
          plt.show()
```
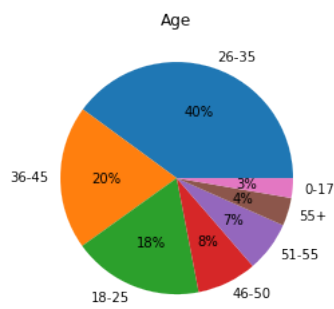


## Observations

- Most of the users are Male
- There are 20 different types of Occupation and Product_Category
- More users belong to B City_Category
- More users are Single as compare to Married

- Product_Category - 1, 5, 8, & 11 have highest purchasing frequency.

In [30]:
```python
data = df['Age'].value_counts(normalize=True)*100
plt.pie(x=data.values, labels=data.index, autopct='%.0f%%')
plt.title("Age")

plt.show()
```
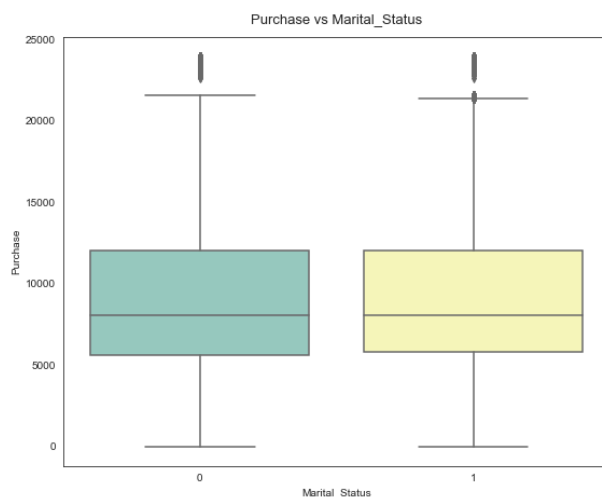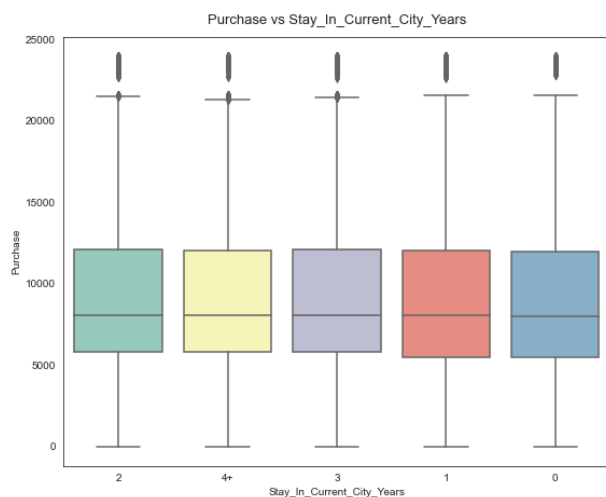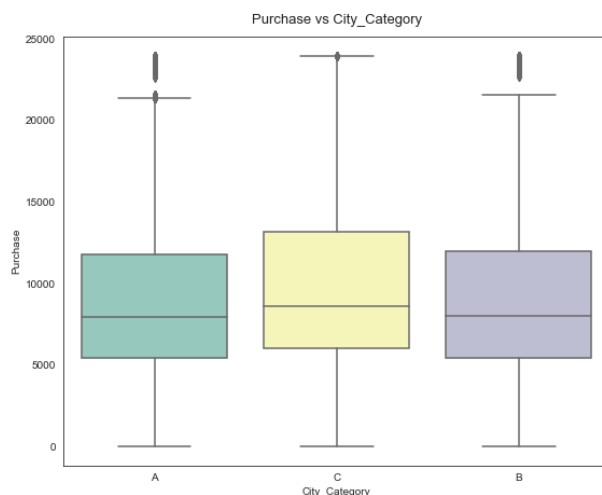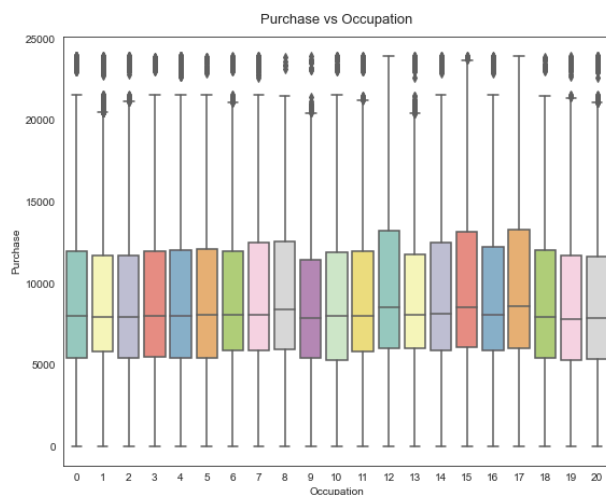
## Bi-variate Analysis

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

fig, axs = plt.subplots(nrows=3, ncols=2, figsize=(20, 16))
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=attrs[count], ax=axs[row, col], palette='Set3')
        axs[row,col].set_title(f"Purchase vs {attrs[count]}", pad=12, fontsize=13)
        count += 1
plt.show()

plt.figure(figsize=(20, 8))
sns.boxplot(data=df, y='Purchase', x=attrs[-1], palette='Set3')
plt.show()
```

```
In [34]: 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', palette='Set3', ax=axs[0,0])
         sns.boxplot(data=df, y='Purchase', x='Gender', hue='City_Category', palette='Set3', ax=axs[0,1])

         sns.boxplot(data=df, y='Purchase', x='Gender', hue='Marital_Status', palette='Set3', ax=axs[1,0])
         sns.boxplot(data=df, y='Purchase', x='Gender', hue='Stay_In_Current_City_Years', palette='Set3', ax=axs[1,1])
         axs[1,1].legend(loc='upper left')

         plt.show()
```



```
In [35]: df.head(10)
```

Out[35]:

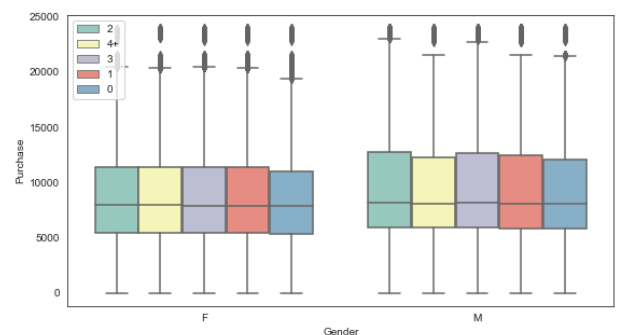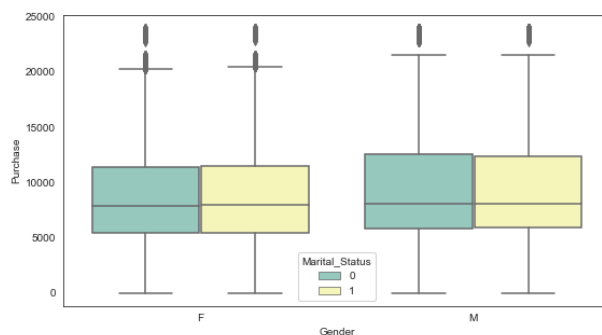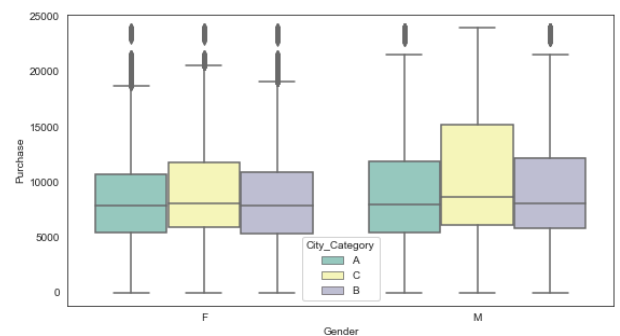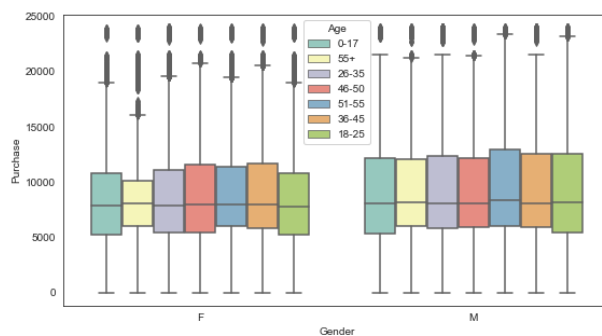|   | User_ID | Product_ID | Gender | Age | Occupation | City_Category | Stay_In_Current_City_Years | Marital_Status | Product_Category | Purchase |
|---|---------|------------|--------|-----|------------|---------------|----------------------------|----------------|------------------|----------|
| 0 | 1000001 | P00069042  | F | 0-17  | 10 | A | 2  | 0 | 3  | 8370  |
| 1 | 1000001 | P00248942  | F | 0-17  | 10 | A | 2  | 0 | 1  | 15200 |
| 2 | 1000001 | P00087842  | F | 0-17  | 10 | A | 2  | 0 | 12 | 1422  |
| 3 | 1000001 | P00085442  | F | 0-17  | 10 | A | 2  | 0 | 12 | 1057  |
| 4 | 1000002 | P00285442  | M | 55+   | 16 | C | 4+ | 0 | 8  | 7969  |
| 5 | 1000003 | P00193542  | M | 26-35 | 15 | A | 3  | 0 | 1  | 15227 |
| 6 | 1000004 | P00184942  | M | 46-50 | 7  | B | 2  | 1 | 1  | 19215 |
| 7 | 1000004 | P00346142  | M | 46-50 | 7  | B | 2  | 1 | 1  | 15854 |
| 8 | 1000004 | P0097242   | M | 46-50 | 7  | B | 2  | 1 | 1  | 15686 |
| 9 | 1000005 | P00274942  | M | 26-35 | 20 | A | 1  | 1 | 8  | 7871  |

**Average amount spend per customer for Male and Female**

In [36]:
```python
amt_df = df.groupby(['User_ID', 'Gender'])[['Purchase']].sum()
amt_df = amt_df.reset_index()
amt_df
```

Out[36]:

|  | 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 [38]:
```python
# Gender wise value counts in avg_amt_df
amt_df['Gender'].value_counts()
```

Out[38]:
```
M    4225
F    1666
Name: Gender, dtype: int64
```

In [39]:
```python
# histogram of average amount spend for each customer - Male & Female
amt_df[amt_df['Gender']=='M']['Purchase'].hist(bins=35)
plt.show()

amt_df[amt_df['Gender']=='F']['Purchase'].hist(bins=35)
plt.show()
```





In [40]:
```python
male_avg = amt_df[amt_df['Gender']=='M']['Purchase'].mean()
female_avg = amt_df[amt_df['Gender']=='F']['Purchase'].mean()

print("Average amount spend by Male customers: {:.2f}".format(male_avg))
print("Average amount spend by Female customers: {:.2f}".format(female_avg))
```

```
Average amount spend by Male customers: 925344.40
Average amount spend by Female customers: 712024.39
```

## Observation

- Male customers spend more money than female customers

```python
In [43]: male_df = amt_df[amt_df['Gender']=='M']
         female_df = amt_df[amt_df['Gender']=='F']
         female_df
```

Out[43]:

|      | User_ID | Gender | Purchase |
|------|---------|--------|----------|
| 0    | 1000001 | F      | 334093   |
| 5    | 1000006 | F      | 379930   |
| 9    | 1000010 | F      | 2169510  |
| 10   | 1000011 | F      | 557023   |
| 15   | 1000016 | F      | 150490   |
| ...  | ...     | ...    | ...      |
| 5885 | 1006035 | F      | 956645   |
| 5886 | 1006036 | F      | 4116058  |
| 5887 | 1006037 | F      | 1119538  |
| 5888 | 1006038 | F      | 90034    |
| 5889 | 1006039 | F      | 590319   |

1666 rows × 3 columns

```python
In [44]: male_df
```

Out[44]:

|      | User_ID | Gender | Purchase |
|------|---------|--------|----------|
| 1    | 1000002 | M      | 810472   |
| 2    | 1000003 | M      | 341635   |
| 3    | 1000004 | M      | 206468   |
| 4    | 1000005 | M      | 821001   |
| 6    | 1000007 | M      | 234668   |
| ...  | ...     | ...    | ...      |
| 5880 | 1006030 | M      | 737361   |
| 5882 | 1006032 | M      | 517261   |
| 5883 | 1006033 | M      | 501843   |
| 5884 | 1006034 | M      | 197086   |
| 5890 | 1006040 | M      | 1653299  |

4225 rows × 3 columns

```python
In [45]: genders = ["M", "F"]

         male_sample_size = 3000
         female_sample_size = 1500
         num_repitions = 1000
         male_means = []
         female_means = []

         for _ in range(num_repitions):
             male_mean = male_df.sample(male_sample_size, replace=True)['Purchase'].mean()
             female_mean = female_df.sample(female_sample_size, replace=True)['Purchase'].mean()

             male_means.append(male_mean)
             female_means.append(female_mean)
```
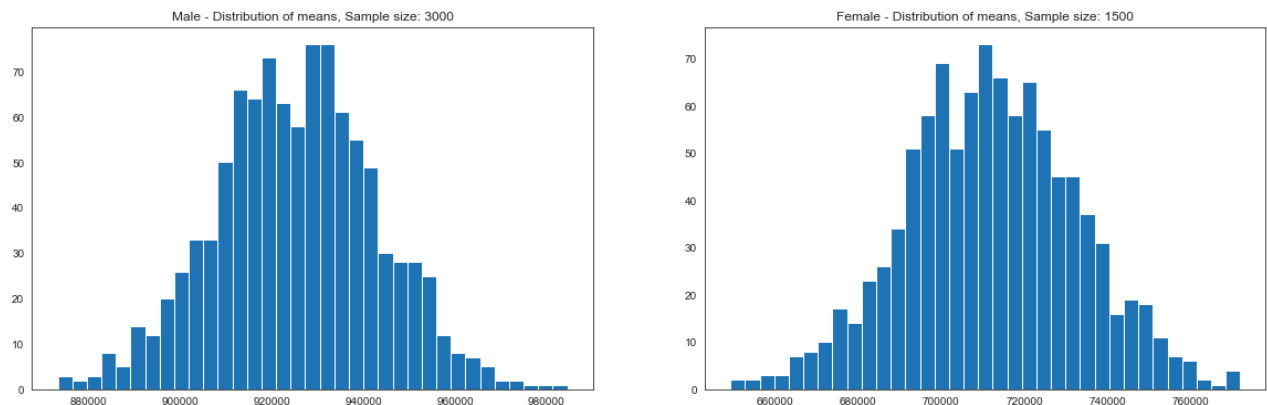
In [46]:
```python
fig, axis = plt.subplots(nrows=1, ncols=2, figsize=(20, 6))

axis[0].hist(male_means, bins=35)
axis[1].hist(female_means, bins=35)
axis[0].set_title("Male - Distribution of means, Sample size: 3000")
axis[1].set_title("Female - Distribution of means, Sample size: 1500")

plt.show()
```



In [47]:
```python
print("Population mean - Mean of sample means of amount spend for Male: {:.2f}".format(np.mean(male_means)))
print("Population mean - Mean of sample means of amount spend for Female: {:.2f}".format(np.mean(female_means)))

print("\nMale - Sample mean: {:.2f} Sample std: {:.2f}".format(male_df['Purchase'].mean(), male_df['Purchase'].std()))
print("Female - Sample mean: {:.2f} Sample std: {:.2f}".format(female_df['Purchase'].mean(), female_df['Purchase'].std()))
```

```
Population mean - Mean of sample means of amount spend for Male: 925268.49
Population mean - Mean of sample means of amount spend for Female: 712000.55

Male - Sample mean: 925344.40 Sample std: 985830.10
Female - Sample mean: 712024.39 Sample std: 807370.73
```

## Observation

Now using the Central Limit Theorem for the population we can say that:

- Average amount spend by male customers is 9,26,341.86
- Average amount spend by female customers is 7,11,704.09

In [48]:
```python
amt_df
```

Out[48]:

|      | 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 [49]:
```python
amt_df = df.groupby(['User_ID', 'Marital_Status'])[['Purchase']].sum()
amt_df = amt_df.reset_index()
amt_df
```

Out[49]:

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

5891 rows × 3 columns

In [50]:
```python
amt_df['Marital_Status'].value_counts()
```

Out[50]:
```
0    3417
1    2474
Name: Marital_Status, dtype: int64
```

In [51]:
```python
marid_samp_size = 3000
unmarid_sample_size = 2000
num_repitions = 1000
marid_means = []
unmarid_means = []

for _ in range(num_repitions):
    marid_mean = amt_df[amt_df['Marital_Status']==1].sample(marid_samp_size, replace=True)['Purchase'].mean()
    unmarid_mean = amt_df[amt_df['Marital_Status']==0].sample(unmarid_sample_size, replace=True)['Purchase'].mean()

    marid_means.append(marid_mean)
    unmarid_means.append(unmarid_mean)


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

axis[0].hist(marid_means, bins=35)
axis[1].hist(unmarid_means, bins=35)
axis[0].set_title("Married - Distribution of means, Sample size: 3000")
axis[1].set_title("Unmarried - Distribution of means, Sample size: 2000")

plt.show()

print("Population mean - Mean of sample means of amount spend for Married: {:.2f}".format(np.mean(marid_means)))
print("Population mean - Mean of sample means of amount spend for Unmarried: {:.2f}".format(np.mean(unmarid_means)))

print("\nMarried - Sample mean: {:.2f} Sample std: {:.2f}".format(amt_df[amt_df['Marital_Status']==1]['Purchase'].mean(), amt
print("Unmarried - Sample mean: {:.2f} Sample std: {:.2f}".format(amt_df[amt_df['Marital_Status']==0]['Purchase'].mean(), amt
```
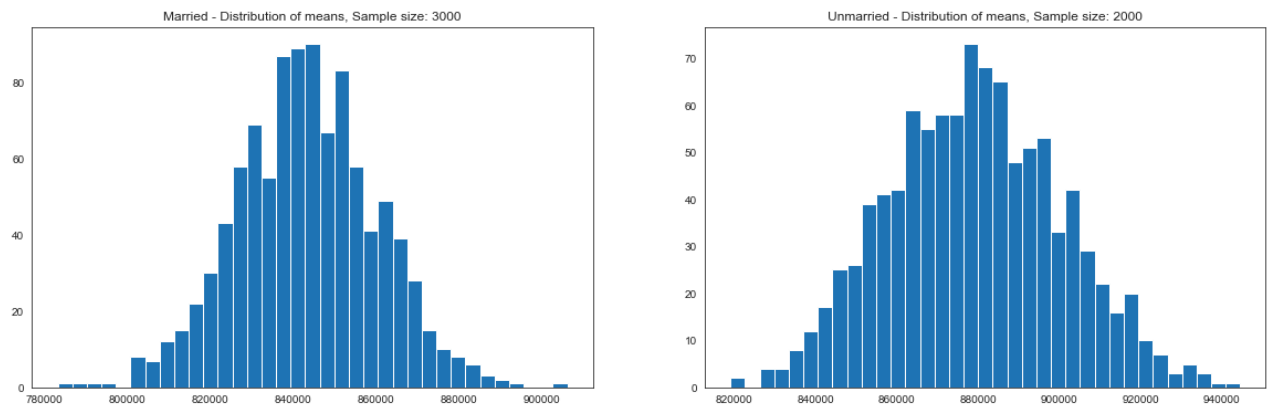


```
Population mean - Mean of sample means of amount spend for Married: 843148.70
Population mean - Mean of sample means of amount spend for Unmarried: 879103.71

Married - Sample mean: 843526.80 Sample std: 935352.12
Unmarried - Sample mean: 880575.78 Sample std: 949436.25
```

In [52]:
```python
amt_df = df.groupby(['User_ID', 'Age'])[['Purchase']].sum()
amt_df = amt_df.reset_index()
amt_df
```

Out[52]:

|  | User_ID | Age | Purchase |
|---|---|---|---|
| 0 | 1000001 | 0-17 | 334093 |
| 1 | 1000002 | 55+ | 810472 |
| 2 | 1000003 | 26-35 | 341635 |
| 3 | 1000004 | 46-50 | 206468 |
| 4 | 1000005 | 26-35 | 821001 |
| ... | ... | ... | ... |
| 5886 | 1006036 | 26-35 | 4116058 |
| 5887 | 1006037 | 46-50 | 1119538 |
| 5888 | 1006038 | 55+ | 90034 |
| 5889 | 1006039 | 46-50 | 590319 |
| 5890 | 1006040 | 26-35 | 1653299 |

5891 rows × 3 columns

In [53]:
```python
amt_df['Age'].value_counts()
```

Out[53]:
```
26-35    2053
36-45    1167
18-25    1069
46-50     531
51-55     481
55+       372
0-17      218
Name: Age, dtype: int64
```

In [54]:
```python
sample_size = 200
num_repitions = 1000

all_means = {}

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

for age_interval in age_intervals:
    for _ in range(num_repitions):
        mean = amt_df[amt_df['Age']==age_interval].sample(sample_size, replace=True)['Purchase'].mean()
        all_means[age_interval].append(mean)
```

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

    new_df = amt_df[amt_df['Age']==val]

    margin_of_error_clt = 1.96*new_df['Purchase'].std()/np.sqrt(len(new_df))
    sample_mean = new_df['Purchase'].mean()
    lower_lim = sample_mean - margin_of_error_clt
    upper_lim = sample_mean + margin_of_error_clt

    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)
```

## Observations

### Confidence Interval by Gender

Now using the Central Limit Theorem for the population:

- Average amount spend by male customers is 9,26,341.86
- Average amount spend by female customers is 7,11,704.09

Now we can infer about the population that, 95% of the times:

- Average amount spend by male customer will lie in between: (895617.83, 955070.97)
- Average amount spend by female customer will lie in between: (673254.77, 750794.02)

### Confidence Interval by Marital_Status

- Married confidence interval of means: (806668.83, 880384.76)
- Unmarried confidence interval of means: (848741.18, 912410.38)

**Confidence Interval by Age**

- 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)

## Recommendations

1. Men spent more money than women, So company should focus on retaining the male customers and getting more male customers.
2. Product_Category - 1, 5, 8, & 11 have highest purchasing frequency. it means these are the products in these categories are liked more by customers. Company can focus on selling more of these products or selling more of the products which are purchased less.
3. Unmarried customers spend more money than married customers, So company should focus on acquisition of Unmarried customers.
4. Customers in the age 18-45 spend more money than the others, So company should focus on acquisition of customers who are in the age 18-45
5. Male customers living in City_Category C spend more money than other male customers living in B or C, Selling more products in the City_Category C will help the company increase the revenue.

In [ ]: