




```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
import warnings
warnings.filterwarnings('ignore')
```

```
df=pd.read_csv('/content/Mall_Customers.csv')
```

```
df.head()
```



	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40





Next steps:



[Generate code with df](#)
[View recommended plots](#)
[New interactive sheet](#)

Univariate Analysis

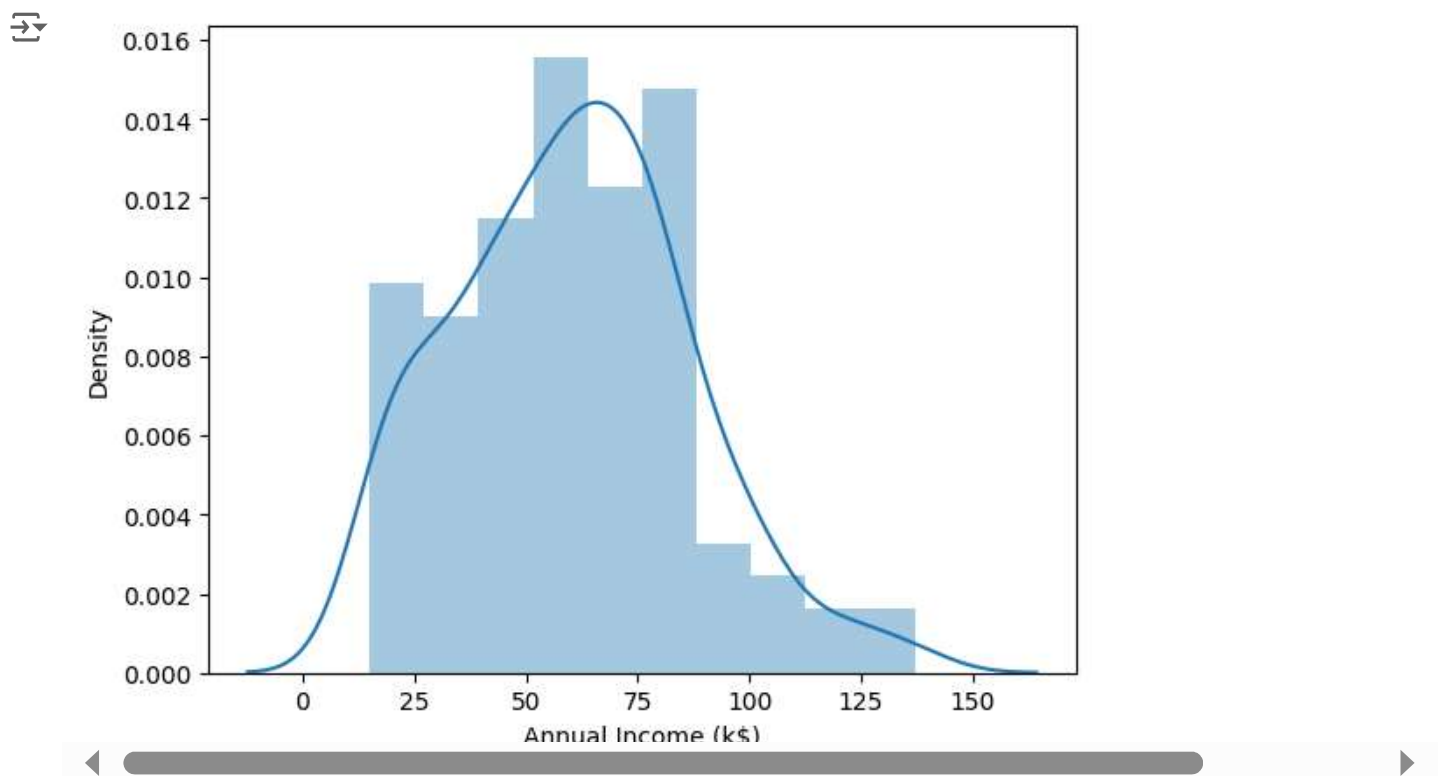
```
df.describe()
```



	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000	200.000000
mean	100.500000	38.850000	60.560000	50.200000
std	57.879185	13.969007	26.264721	25.823522
min	1.000000	18.000000	15.000000	1.000000
25%	50.750000	28.750000	41.500000	34.750000
50%	100.500000	36.000000	61.500000	50.000000
75%	150.250000	49.000000	78.000000	73.000000
max	200.000000	70.000000	137.000000	99.000000

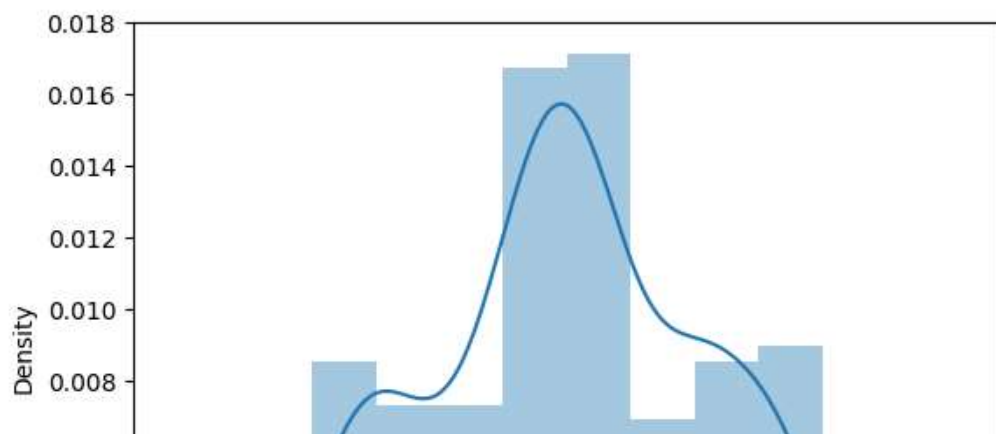
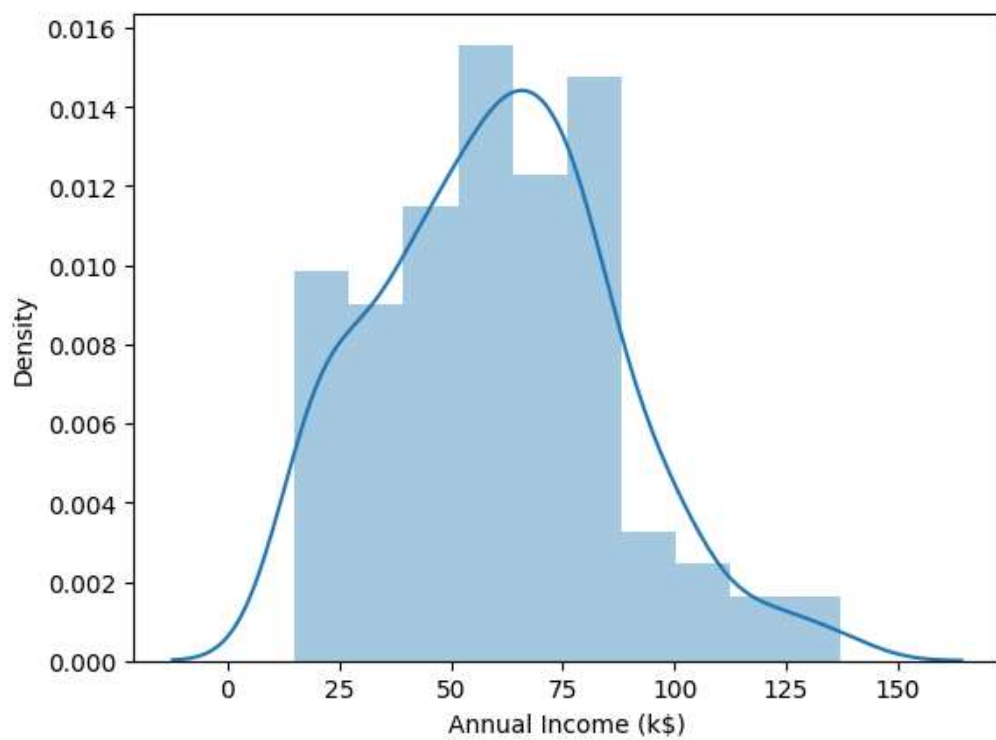
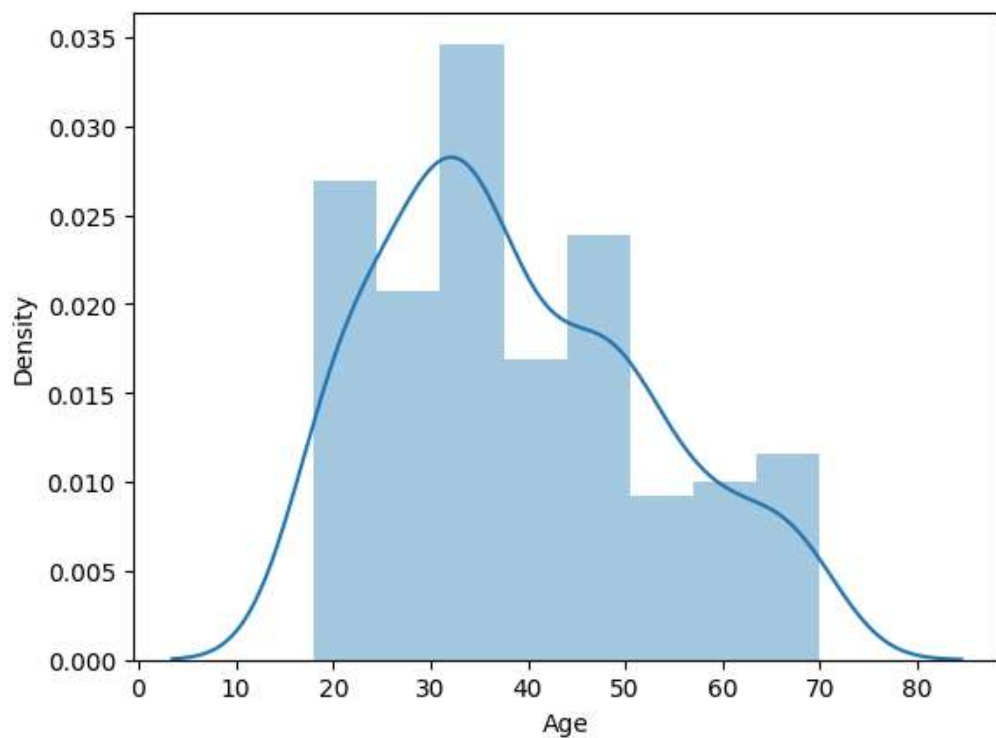
```
sns.distplot(df['Annual Income (k$)']);
```

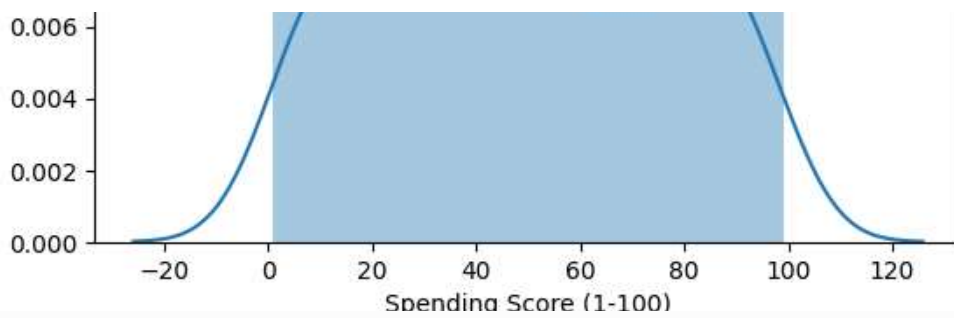


```
df.columns
```

```
Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',  
      'Spending Score (1-100)'],  
      dtype='object')
```

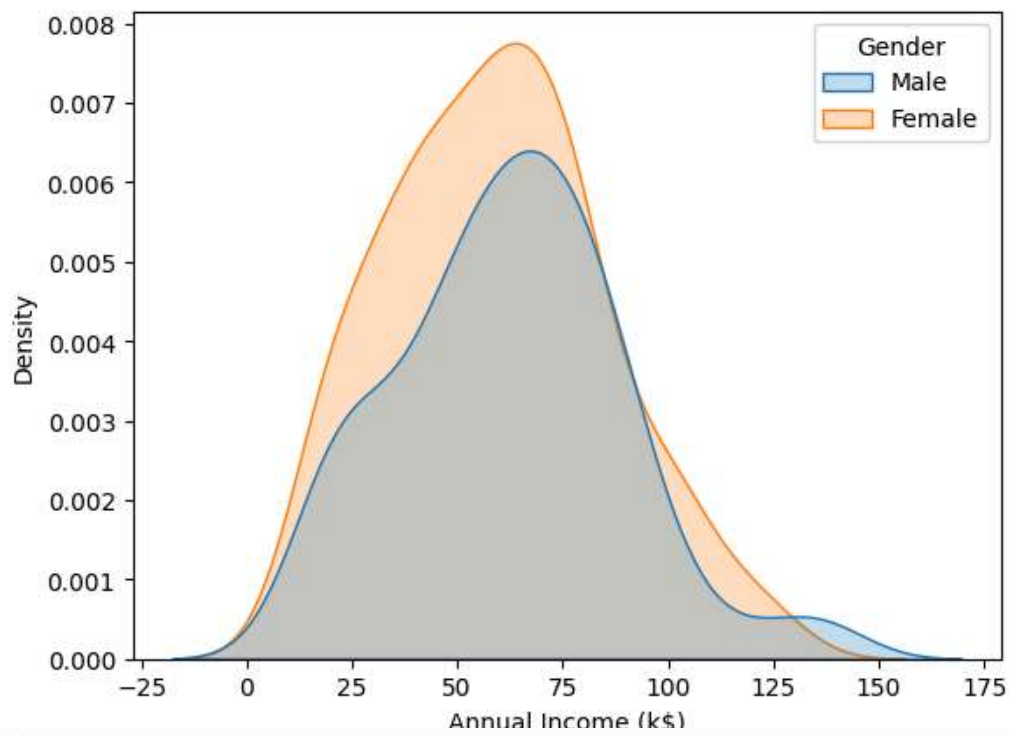
```
columns = ['Age', 'Annual Income (k$)', 'Spending Score (1-100)']  
for i in columns:  
    plt.figure()  
    sns.distplot(df[i])
```



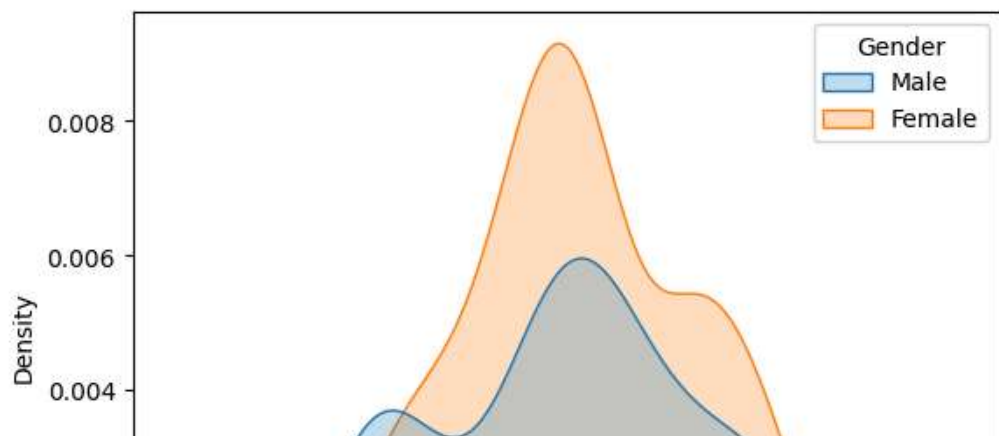
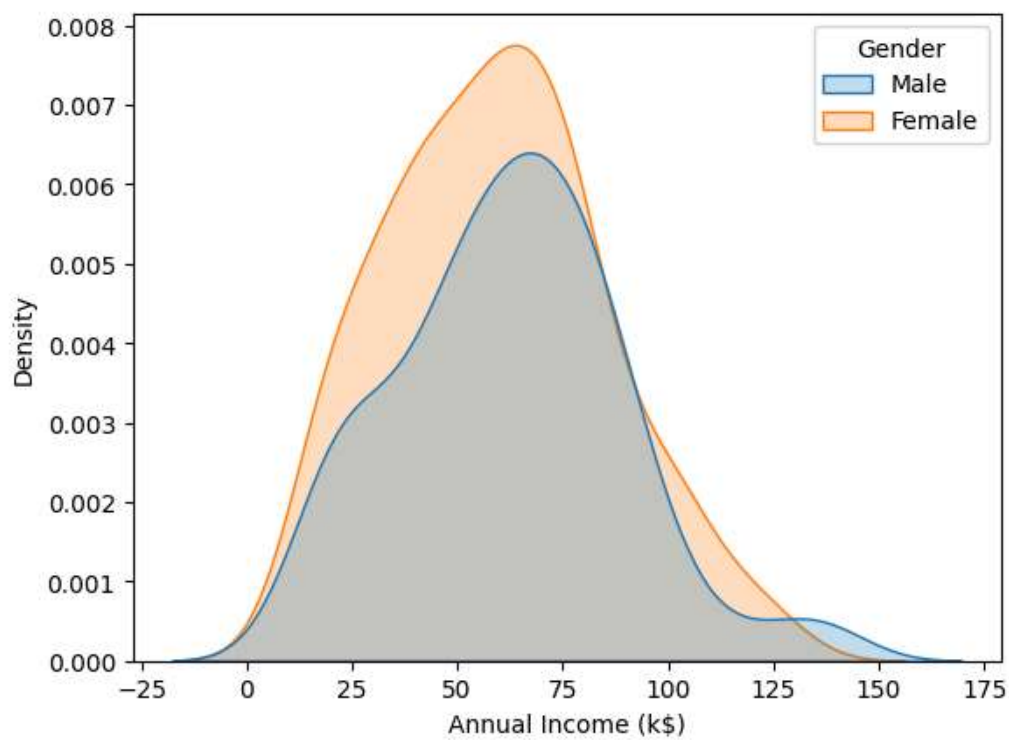
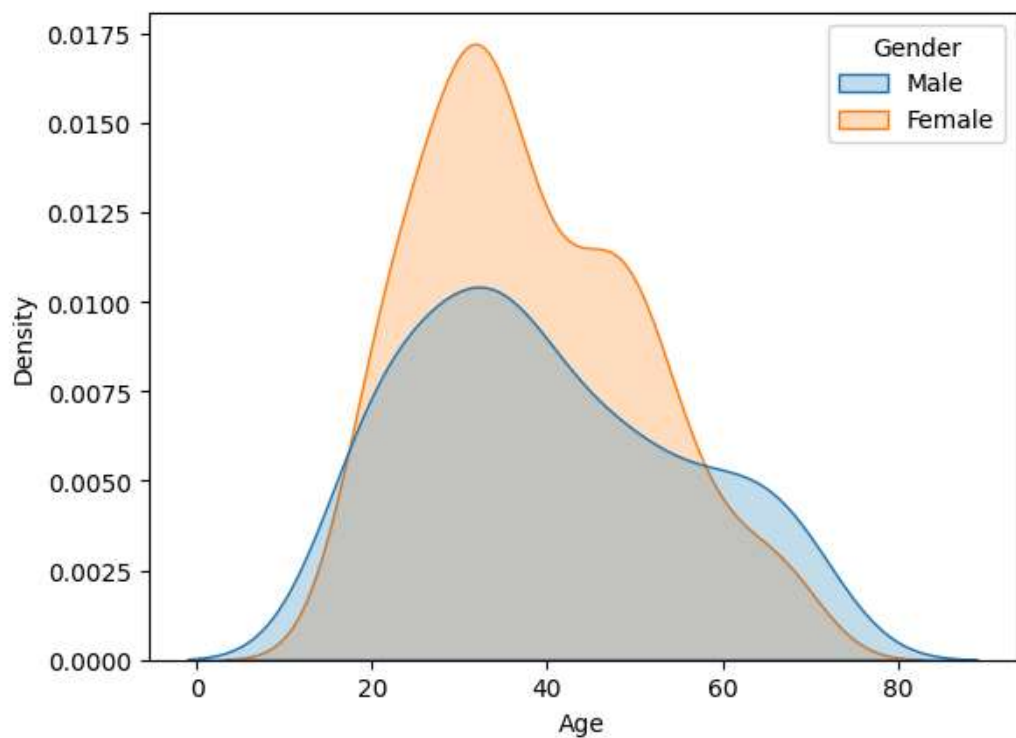


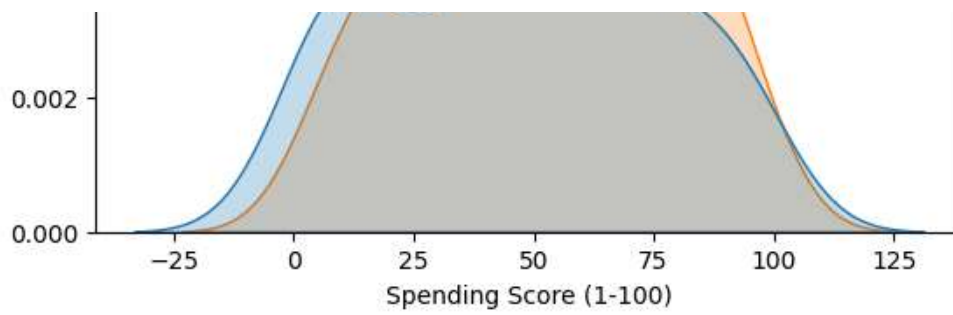
```
sns.kdeplot(data=df, x='Annual Income (k$)', hue='Gender', shade=True)
```

<Axes: xlabel='Annual Income (k\$)', ylabel='Density'>

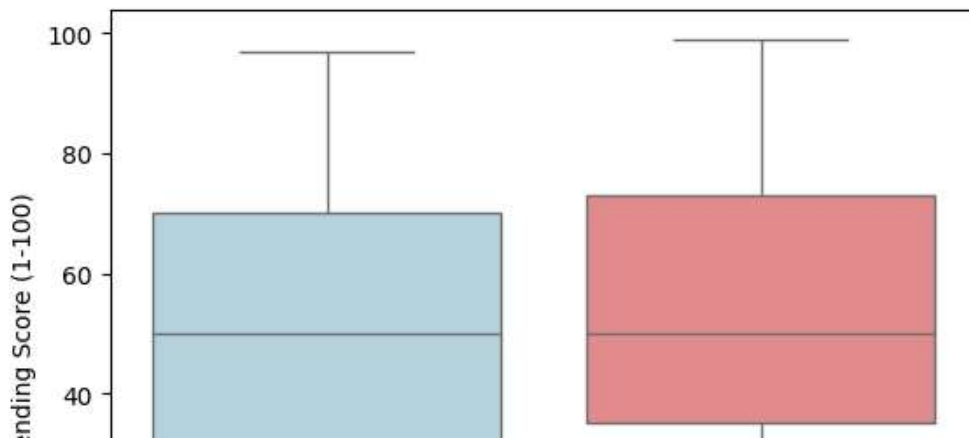
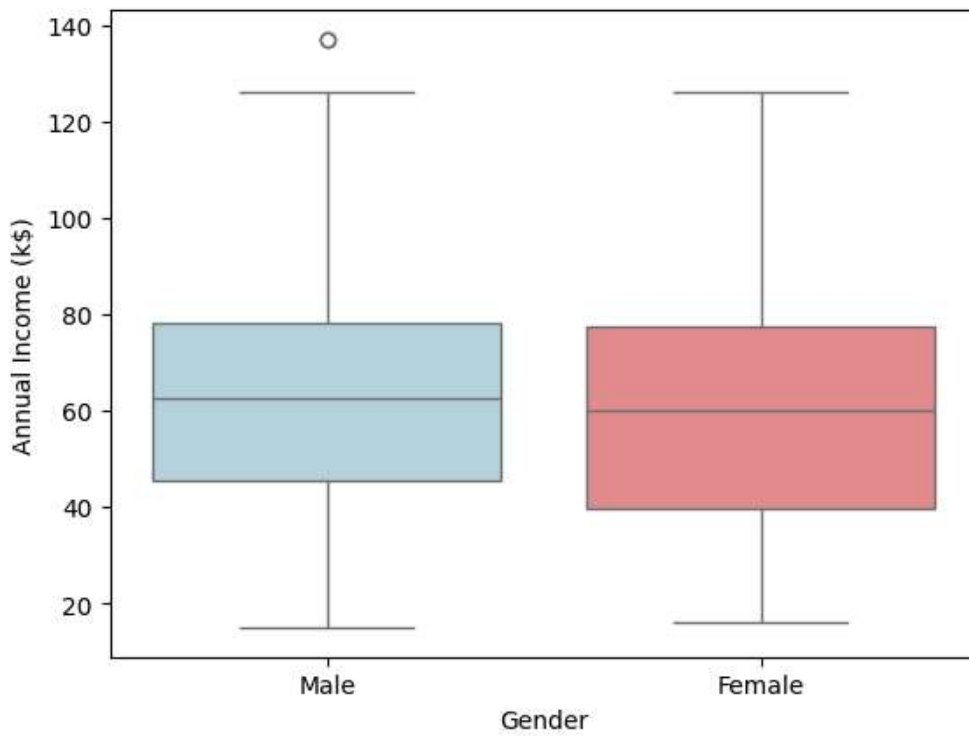
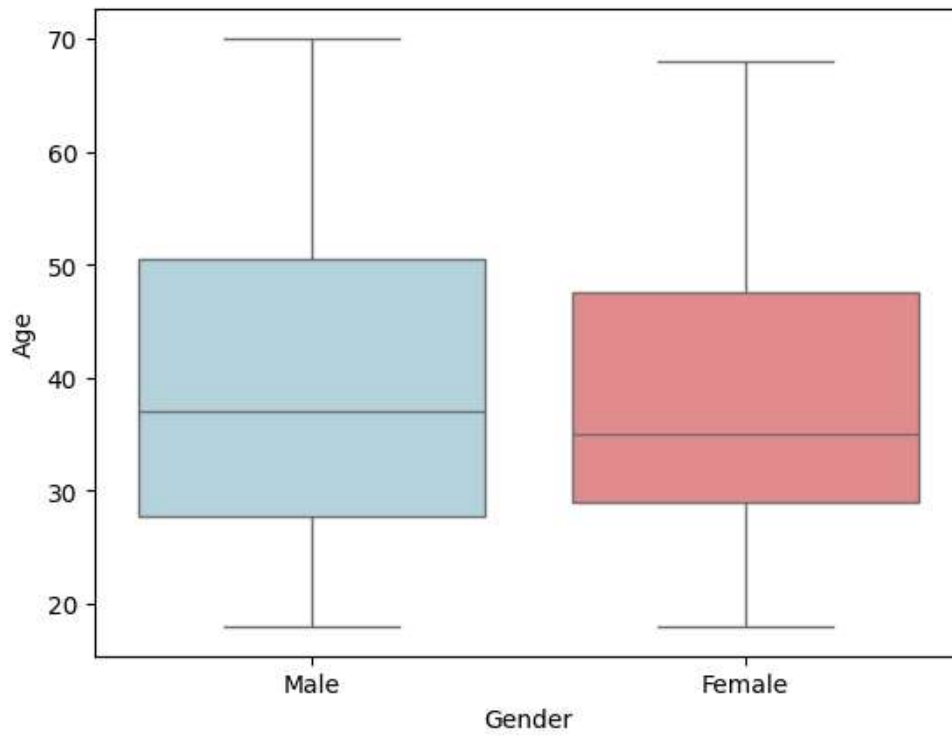


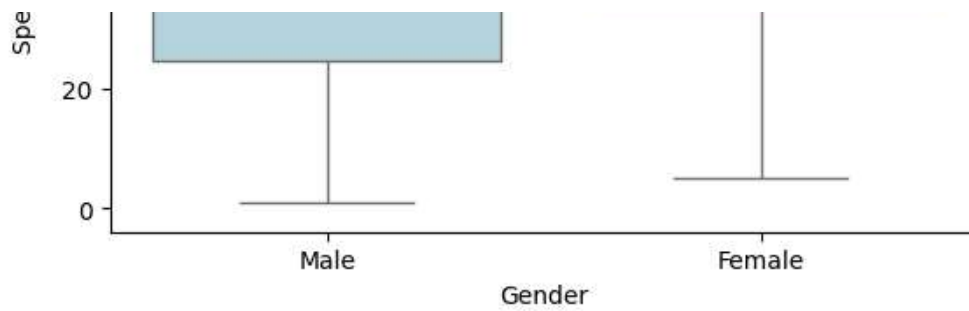
```
columns = ['Age', 'Annual Income (k$)', 'Spending Score (1-100)']
for i in columns:
    plt.figure()
    sns.kdeplot(data=df, x=i, hue='Gender', shade=True)
```





```
columns = ['Age', 'Annual Income (k$)', 'Spending Score (1-100)']  
for i in columns:  
    plt.figure()  
    sns.boxplot(data=df, x='Gender', y=df[i], palette=['lightblue', 'lightcoral'])
```





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



proportion	
Gender	
Female	0.56
Male	0.44

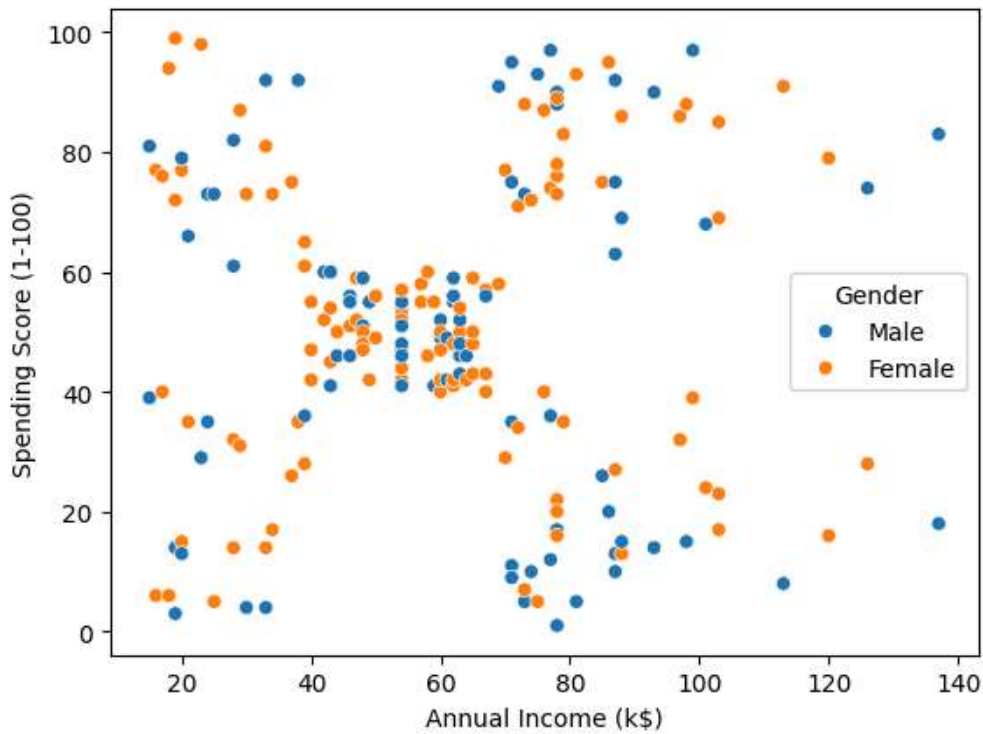
dtype: float64

✓ Bivariate Analysis

```
sns.scatterplot(data=df, x='Annual Income (k$)', y='Spending Score (1-100)', hue='Gender')
```



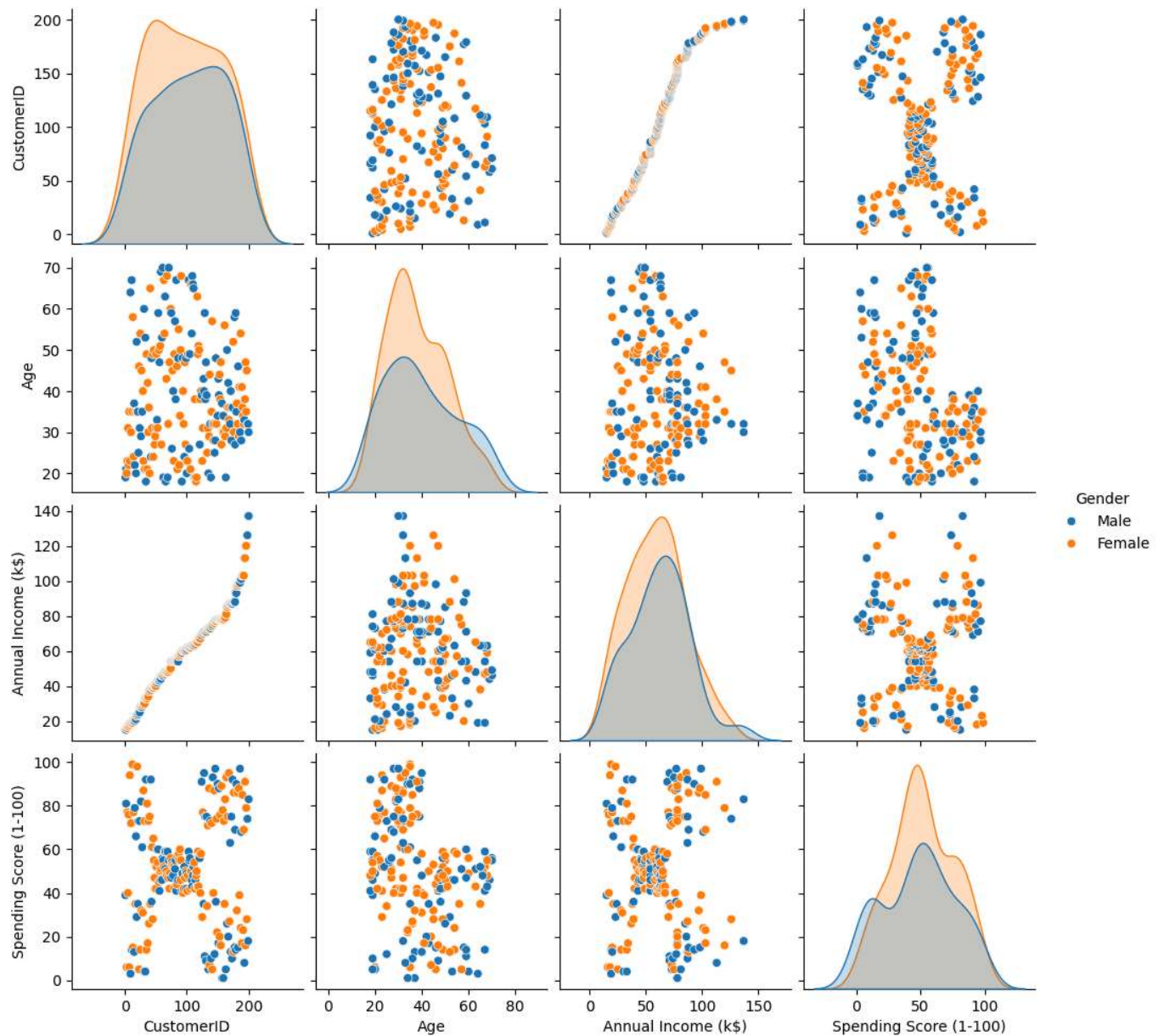
<Axes: xlabel='Annual Income (k\$)', ylabel='Spending Score (1-100)'>



```
#df=df.drop('CustomerID',axis=1)
sns.pairplot(df,hue='Gender')
```




```
>>> <seaborn.axisgrid.PairGrid at 0x788f40ebe830>
```

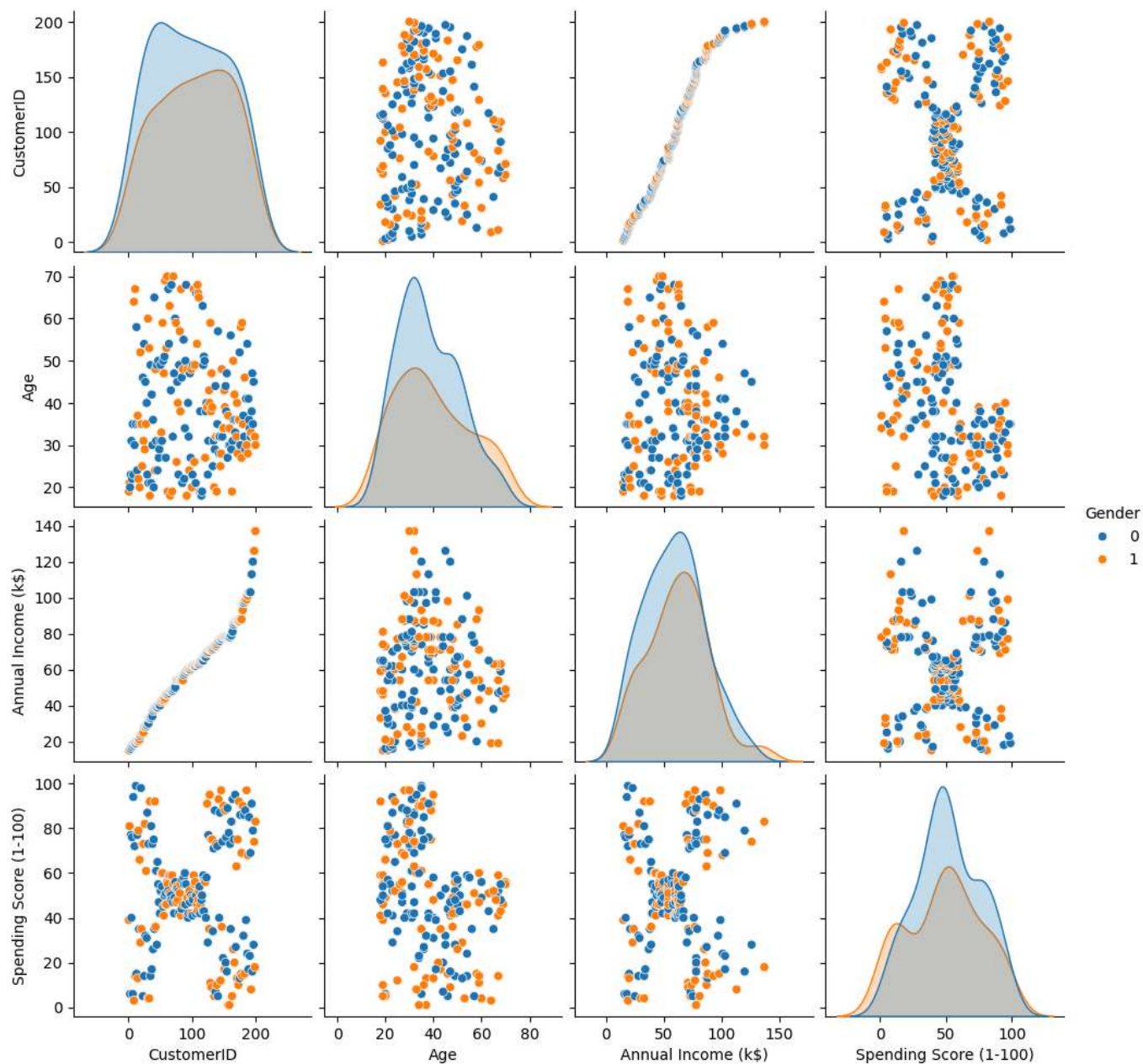


```
df.groupby(['Gender'])[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']].mean()
```

	Age	Annual Income (k\$)	Spending Score (1-100)
Gender			
Female	38.098214	59.250000	51.526786
Male	39.806818	62.227273	48.511364

```
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['Gender'] = le.fit_transform(df['Gender'])
sns.pairplot(df, hue='Gender')
```

 <seaborn.axisgrid.PairGrid at 0x788f3e3658d0>



df.corr()



	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
CustomerID	1.000000	0.057400	-0.026763	0.977548	0.013835
Gender	0.057400	1.000000	0.060867	0.056410	-0.058109
Age	-0.026763	0.060867	1.000000	-0.012398	-0.327227
Annual Income (k\$)	0.977548	0.056410	-0.012398	1.000000	0.009903
Spending Score (1-	0.013835	-0.058109	-0.327227	0.009903	1.000000



 <Axes: >



```
clustering1 = KMeans(n_clusters=3)
```

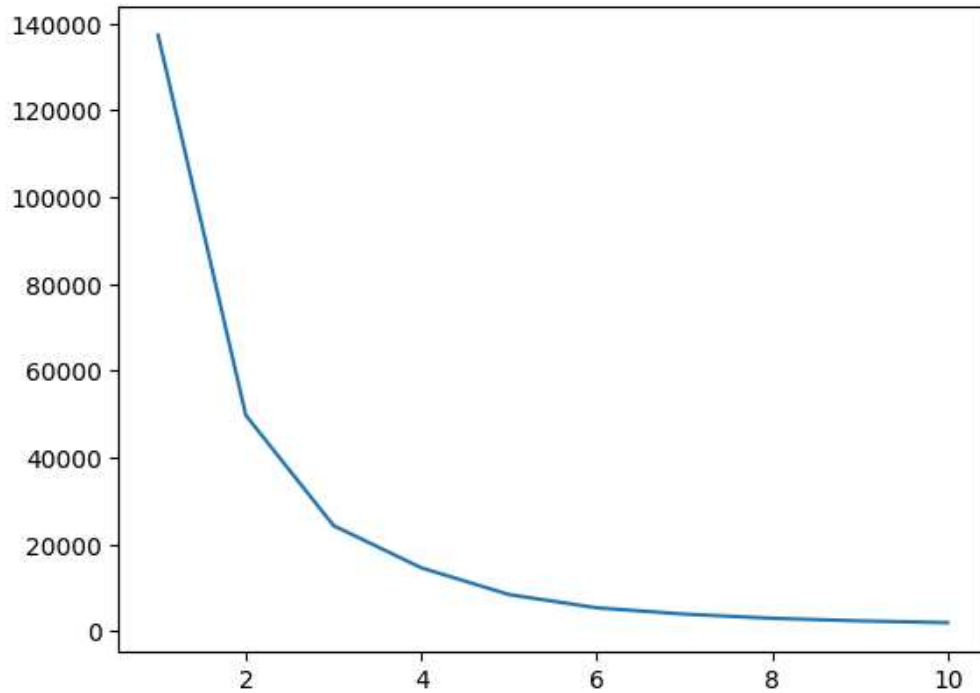
▼ KMeans ⓘ ?
KMeans(n_clusters=3)

```
 array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,  
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,  
       2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 0, 0, 0, 0, 0, 0, 0, 0,  
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,  
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,  
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,  
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1])
```



```
plt.plot(range(1,11),intertia_scores)
```

```
[<matplotlib.lines.Line2D at 0x788f3d916c80>]
```



```
df.columns
```

```
Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',  
      'Spending Score (1-100)', 'Income Cluster'],  
      dtype='object')
```

```
df.groupby('Income Cluster')[['Age', 'Annual Income (k$)',  
                              'Spending Score (1-100)']].mean()
```

	Age	Annual Income (k\$)	Spending Score (1-100)
Income Cluster			
0	41.279070	60.906977	50.337209
1	36.910714	92.142857	50.517857
2	37.120690	29.551724	49.689655

▼ Bivariate Clustering

```
clustering2 = KMeans(n_clusters=5)
clustering2.fit(df[['Annual Income (k$)', 'Spending Score (1-100)']])
df['Spending and Income Cluster'] =clustering2.labels_
df.head()
```



	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)	Income Cluster	Spending and Income Cluster
0	1	1	19	15	39	2	0
1	2	1	21	15	81	2	1
2	3	0	20	16	6	2	0
3	4	0	23	16	77	2	1
4	5	0	31	17	40	2	0



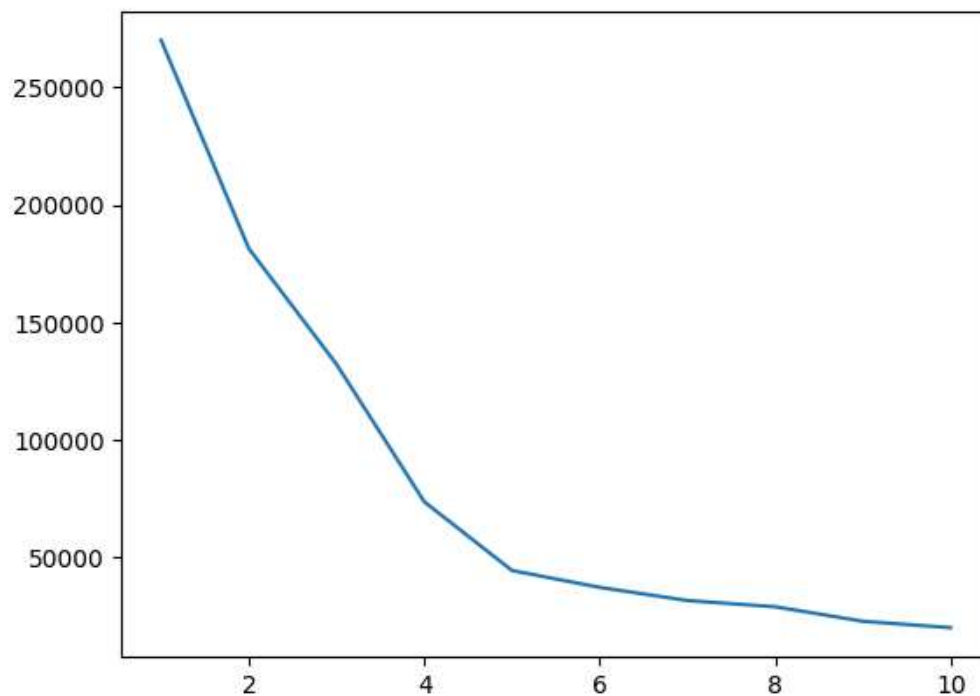
Next steps:

[Generate code with df](#)[View recommended plots](#)[New interactive sheet](#)

```
intertia_scores2=[]
for i in range(1,11):
    kmeans2=KMeans(n_clusters=i)
    kmeans2.fit(df[['Annual Income (k$)', 'Spending Score (1-100)']])
    inertia_scores2.append(kmeans2.inertia_)
plt.plot(range(1,11),intertia_scores2)
```

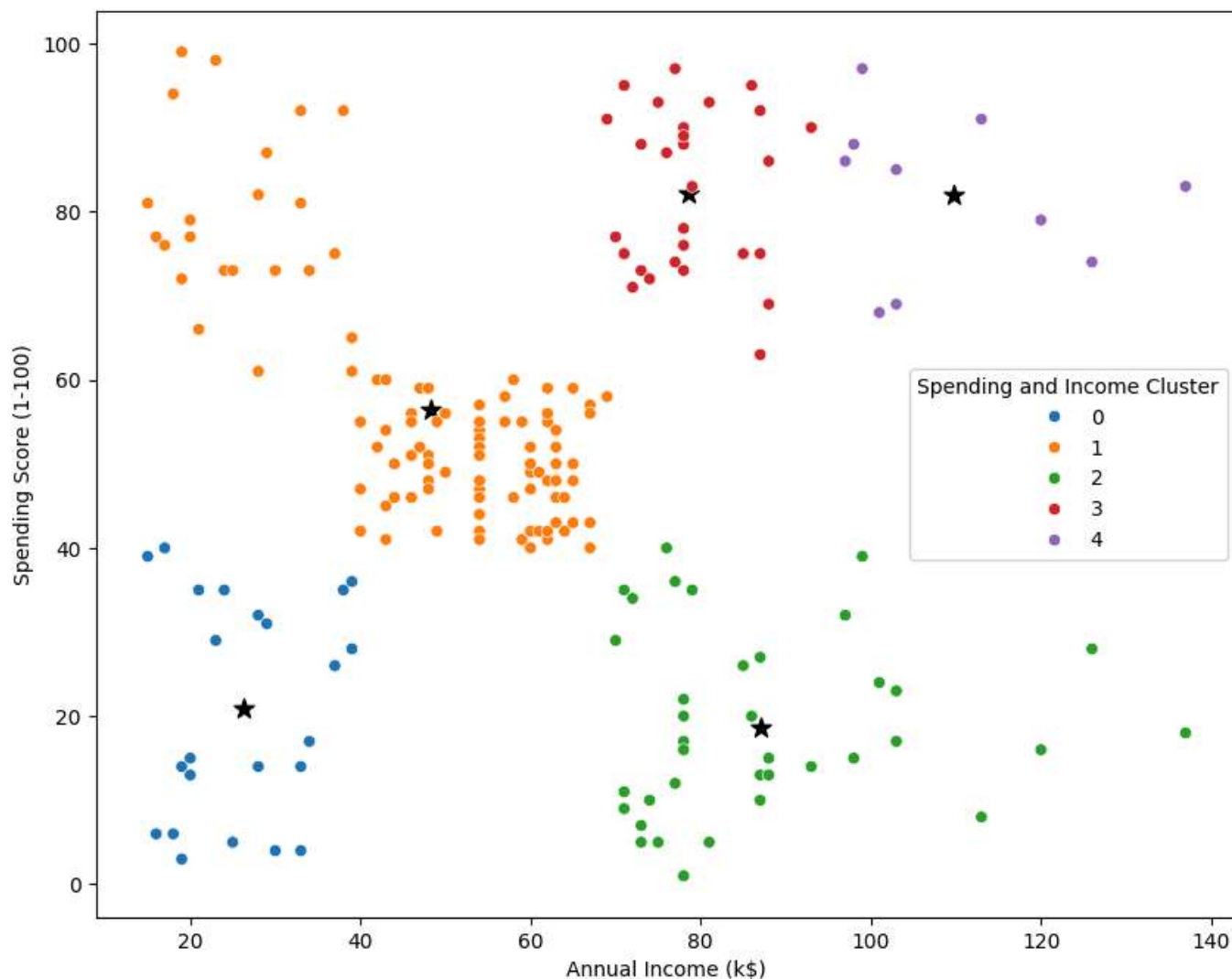


[<matplotlib.lines.Line2D at 0x788f3d7fd810>]



```
centers =pd.DataFrame(clustering2.cluster_centers_)
centers.columns = ['x','y']
```

```
plt.figure(figsize=(10,8))
plt.scatter(x=centers['x'],y=centers['y'],s=100,c='black',marker='*')
sns.scatterplot(data=df, x = 'Annual Income (k$)',y='Spending Score (1-100)',hue='Spending and Income Clust
plt.savefig('clustering_bivaraiate.png')
```



```
pd.crosstab(df['Spending and Income Cluster'],df['Gender'],normalize='index')
```



Spending and Income Cluster	Gender	
	0	1
0	0.608696	0.391304
1	0.590000	0.410000
2	0.473684	0.526316
3	0.517241	0.482759
4	0.600000	0.400000

Generate

ValueError: Cannot subset columns with a tuple with more than one element. Use a list instead.
df.groupby('Spending and Income Cluster')['Age', 'Annual Income (k\$)',




Close

< 1 of 1 >



[Undo Changes](#) [Use code with caution](#)




```
df.groupby('Spending and Income Cluster')[['Age', 'Annual Income (k$)',
'Spending Score (1-100)']].mean()
```




	Age	Annual Income (k\$)	Spending Score (1-100)
Spending and Income Cluster			
0	45.217391	26.304348	20.913043
1	39.200000	48.260000	56.480000
2	40.394737	87.000000	18.631579
3	32.862069	78.551724	82.172414
4	32.200000	109.700000	82.000000





▼ **Multivariate clustering**

```
from sklearn.preprocessing import StandardScaler
scale = StandardScaler()
df.head()
```



	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)	Income Cluster	Spending and Income Cluster
0	1	1	19	15	39	2	0
1	2	1	21	15	81	2	1
2	3	0	20	16	6	2	0
3	4	0	23	16	77	2	1
4	5	0	31	17	40	2	0




Next steps:

[Generate code with df](#)



 [View recommended plots](#)

[New interactive sheet](#)

```
dff = pd.get_dummies(df,drop_first=True)
dff.head()
```



	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)	Income Cluster	Spending and Income Cluster
0	1	1	19	15	39	2	0
1	2	1	21	15	81	2	1
2	3	0	20	16	6	2	0
3	4	0	23	16	77	2	1
4	5	0	31	17	40	2	0




Next steps:

[Generate code with dff](#)

 [View recommended plots](#)



[New interactive sheet](#)


```
dff.columns
```


 `Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',
 'Spending Score (1-100)', 'Income Cluster',
 'Spending and Income Cluster'],
 dtype='object')`

```
dff = dff[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']]
```

```
dff.head()
```

	Age	Annual Income (k\$)	Spending Score (1-100)	
0	19	15	39	
1	21	15	81	
2	20	16	6	
3	23	16	77	
4	31	17	40	

Next steps:

[Generate code with dff](#)

 [View recommended plots](#)

[New interactive sheet](#)

```
dff = scale.fit_transform(dff)
```