In [1]:

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

In [2]:

```python
# Importing the dataset
dataset = pd.read_csv('C:/Users/Jayen/Desktop/study material/2nd year-3sem/machine
learning/programs/Hierarchical-Clustering/Hierarchical_Clustering/Mall_Customers.csv')
dataset=pd.DataFrame(dataset)
dataset.head()
```

Out[2]:

|   | CustomerID | Genre | 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 |

In [3]:

```python
dataset.describe()
```

Out[3]:

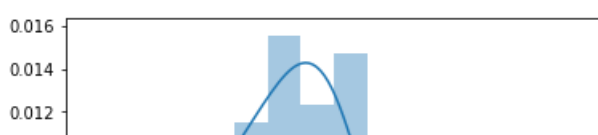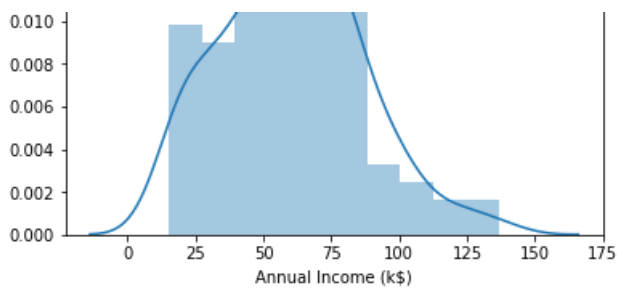|   | CustomerID | Age | Annual Income (k$) | Spending Score (1-100) |
|---|---|---|---|---|
| count | 200.000000 | 200.000000 | 200.000000 | 200.000000 |
| mean | 100.500000 | 38.850000 | 60.560000 | 50.200000 |
| std | 57.879185 | 13.969007 | 26.264721 | 25.823522 |
| min | 1.000000 | 18.000000 | 15.000000 | 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 |

In [4]:

```python
sns.distplot(dataset['Annual Income (k$)'])
```

```
C:\Users\Jayen\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
  warnings.warn("The 'normed' kwarg is deprecated, and has been "
```

Out[4]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x235bccba470>
```

Annual Income (k$)

```
sns.distplot(dataset['Spending Score (1-100)'])
```

```
C:\Users\Jayen\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
  warnings.warn("The 'normed' kwarg is deprecated, and has been "
```

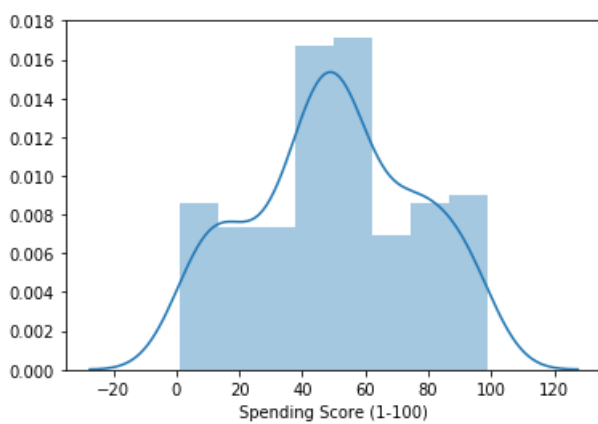Out[5]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x235bd0cf6d8>
```
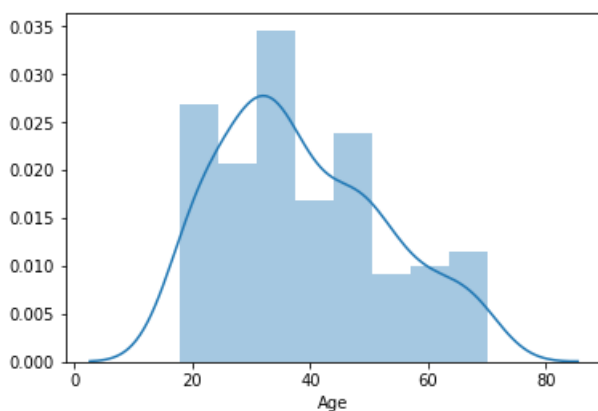


Spending Score (1-100)

In [6]:

```
sns.distplot(dataset['Age'])
```

```
C:\Users\Jayen\Anaconda3\lib\site-packages\matplotlib\axes\_axes.py:6462: UserWarning: The
'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.
  warnings.warn("The 'normed' kwarg is deprecated, and has been "
```
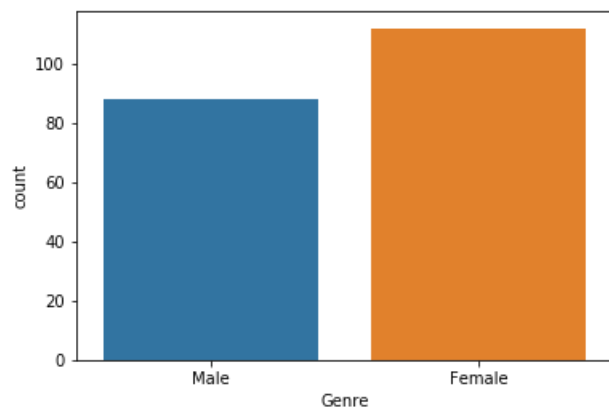
Out[6]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x235be12c5f8>
```



Age

In [7]:

```
sns.countplot(x='Genre',data=dataset)
```

Out[7]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x235be159fd0>
```

```python
cor=dataset.corr()
sns.heatmap(cor,xticklabels=cor.columns.values,yticklabels=cor.columns.values,annot=True)
```

Out[9]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x235be250cf8>
```



In [34]:

```python
X = dataset.iloc[:,4].values.reshape(-1,1)
y = dataset.iloc[:,3].values.reshape(-1,1)
```

In [35]:

```python
# Splitting the dataset into the Training set and Test set
from sklearn.cross_validation import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 0)
```
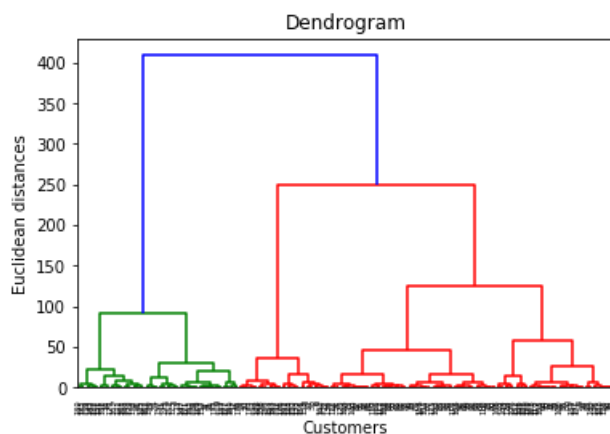
In [36]:

```python
# Feature Scaling
from sklearn.preprocessing import StandardScaler
sc_X = StandardScaler()
X_train = sc_X.fit_transform(X_train)
X_test = sc_X.transform(X_test)
sc_y = StandardScaler()
y_train = sc_y.fit_transform(y_train)
```

```
C:\Users\Jayen\Anaconda3\lib\site-packages\sklearn\utils\validation.py:475: DataConversionWarning:
Data with input dtype int64 was converted to float64 by StandardScaler.
```

In [37]:

```python
# Using the dendrogram to find the optimal number of clusters
import scipy.cluster.hierarchy as sch
dendrogram = sch.dendrogram(sch.linkage(X, method = 'ward'))
plt.title('Dendrogram')
plt.xlabel('Customers')
plt.ylabel('Euclidean distances')
plt.show()
```



In [32]:

```python
# Fitting Hierarchical Clustering to the dataset
from sklearn.cluster import AgglomerativeClustering
hc = AgglomerativeClustering(n_clusters =2, affinity = 'euclidean', linkage = 'ward')
y_hc = hc.fit_predict(X)
```

In [33]:

```python
plt.scatter(X[y_hc == 0, 0], X[y_hc == 0,1], s = 100, c = 'red', label = 'Cluster 1')
plt.scatter(X[y_hc == 1, 0], X[y_hc == 1, 1], s = 100, c = 'blue', label = 'Cluster 2')
plt.scatter(X[y_hc == 2, 0], X[y_hc == 2, 1], s = 100, c = 'green', label = 'Cluster 3')
plt.scatter(X[y_hc == 3, 0], X[y_hc == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4')
plt.scatter(X[y_hc == 4, 0], X[y_hc == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5')
plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```

0

20    40    60    80    100    120    140
Annual Income (k$)

In [65]: