

Descriptive Data Analysis:

Descriptive data analysis primarily involves summarizing and understanding the basic characteristics and features of the dataset without making in-depth inferences or predictions.

1. What is the age distribution in the dataset?
2. What is the distribution of estimated salaries?
3. How many people made a purchase (Purchased=1) and how many did not (Purchased=0)?
4. What is the overall percentage of people who made a purchase?
5. What is the average age and estimated salary of those who made a purchase compared to those who did not?
6. Are there any missing values in the dataset, and if so, in which columns?
7. What is the gender distribution in the dataset (if not included in the provided subset)?
8. What is the relationship between age and estimated salary?

Exploratory Data Analysis:

Exploratory data analysis goes a step further by investigating relationships, patterns, and trends within the data. It often includes visualization and statistical techniques to uncover insights and identify interesting patterns.

1. Is there a correlation between age and the likelihood of making a purchase?
2. Is there a correlation between estimated salary and the likelihood of making a purchase?
3. How does the distribution of estimated salary differ between those who made a purchase and those who did not?
4. Are there any outliers in the dataset, and how do they affect the analysis?

Segmentation and Customer Profiling:

Segmentation and customer profiling involve dividing your data into distinct groups or segments based on certain characteristics and understanding the unique traits and behaviors of these groups.

1. Can you segment the data into different customer groups based on age and estimated salary?
2. What are the characteristics of the customer segments, and how do they differ in terms of purchase behavior?

Feature Importance Analysis:

Feature importance analysis is often part of predictive modeling and understanding which features or variables have the most influence on the target variable (in your case, purchase decisions). It helps identify the key drivers of the outcome. The question related to feature importance is:

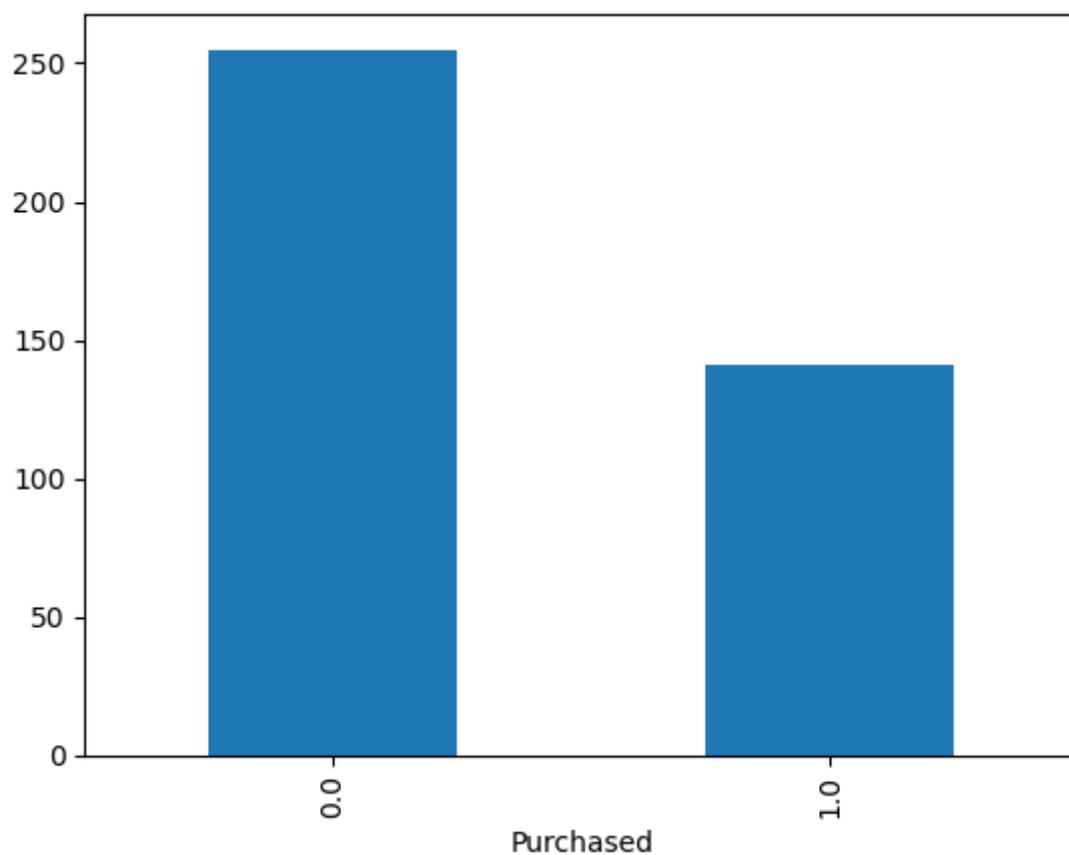
1. Which features (age, estimated salary, etc.) have the most significant impact on purchase decisions?

In [1]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Load your dataset
6 data = pd.read_csv('Social_Network_Ads.csv')
```

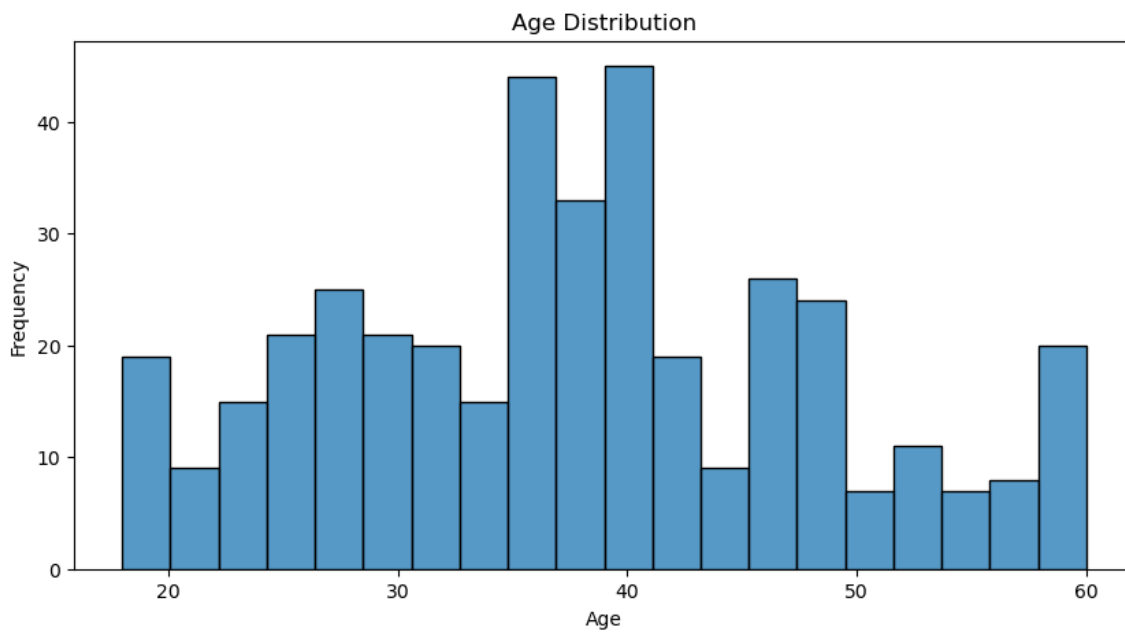
In [12]:

```
1 data["Purchased"].value_counts().plot(kind="bar");
```



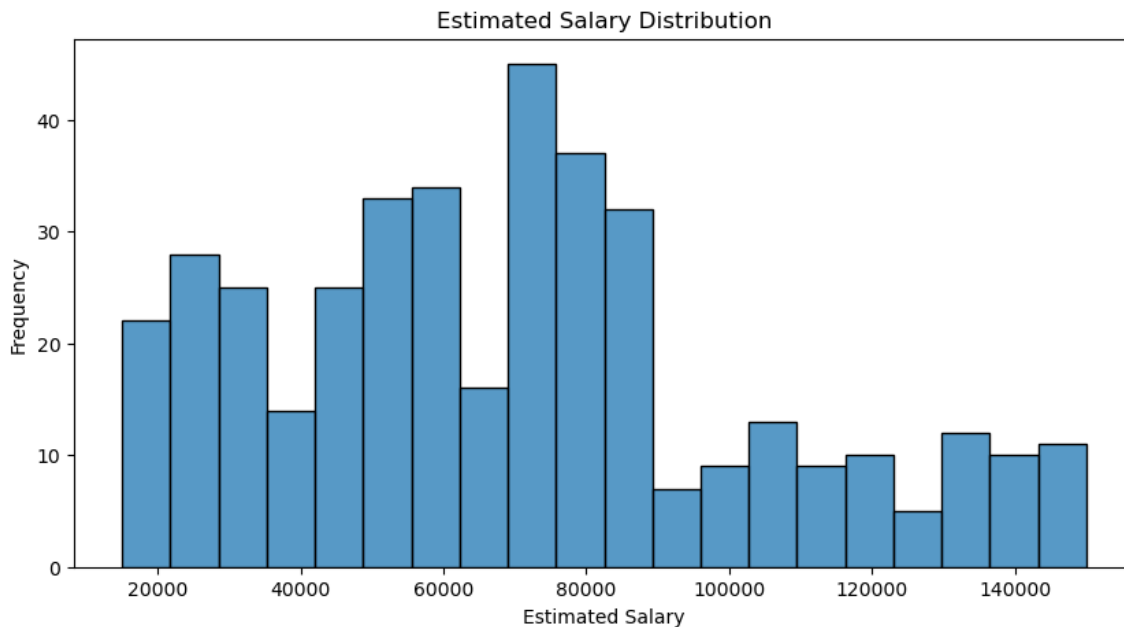
In [3]:

```
1 # Descriptive Questions:
2
3 # What is the age distribution in the dataset?
4 plt.figure(figsize=(10, 5))
5 sns.histplot(data['Age'], bins=20)
6 plt.title('Age Distribution')
7 plt.xlabel('Age')
8 plt.ylabel('Frequency')
9 plt.show()
10
```



In [4]:

```
1 # What is the distribution of estimated salaries?
2 plt.figure(figsize=(10, 5))
3 sns.histplot(data['EstimatedSalary'], bins=20)
4 plt.title('Estimated Salary Distribution')
5 plt.xlabel('Estimated Salary')
6 plt.ylabel('Frequency')
7 plt.show()
```



In [5]:

```
1 # How many people made a purchase (Purchased=1) and how many did not (Purchased=0)?
2 purchase_counts = data['Purchased'].value_counts()
3 print('Purchase Counts:\n', purchase_counts)
4
5
```

Purchase Counts:

Purchased

0.0 255

1.0 141

Name: count, dtype: int64

In [6]:

```
1 # What is the overall percentage of people who made a purchase?
2 percentage_purchase = (purchase_counts[1] / purchase_counts.sum()) * 100
3 print(f'Percentage of people who made a purchase: {percentage_purchase:.2f}%')
4
```

Percentage of people who made a purchase: 35.61%

In [7]:

```
1
2 # What is the average age and estimated salary of those who made a purchase c
3 average_age_purchase = data[data['Purchased'] == 1]['Age'].mean()
4 average_salary_purchase = data[data['Purchased'] == 1]['EstimatedSalary'].mean()
5 average_age_no_purchase = data[data['Purchased'] == 0]['Age'].mean()
6 average_salary_no_purchase = data[data['Purchased'] == 0]['EstimatedSalary'].mean()
7
8 print('Average Age of Purchase:', average_age_purchase)
9 print('Average Estimated Salary of Purchase:', average_salary_purchase)
10 print('Average Age of No Purchase:', average_age_no_purchase)
11 print('Average Estimated Salary of No Purchase:', average_salary_no_purchase)
12
```

Average Age of Purchase: 46.51063829787234

Average Estimated Salary of Purchase: 86338.12949640288

Average Age of No Purchase: 32.77865612648221

Average Estimated Salary of No Purchase: 60444.88188976378

In [8]:

```
1 # Are there any missing values in the dataset, and if so, in which columns?
2 missing_values = data.isnull().sum()
3 print('Missing Values:\n', missing_values)
4
5
```

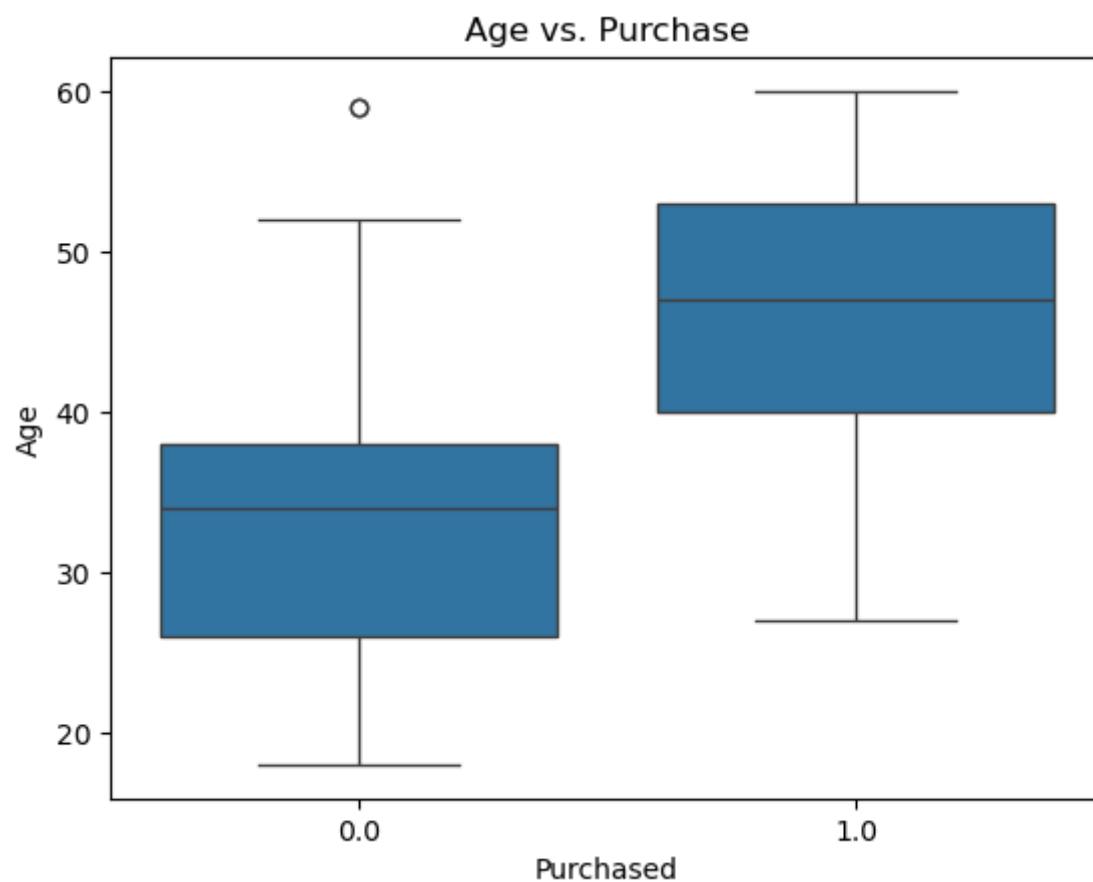
Missing Values:

Age	2
EstimatedSalary	3
Purchased	4

dtype: int64

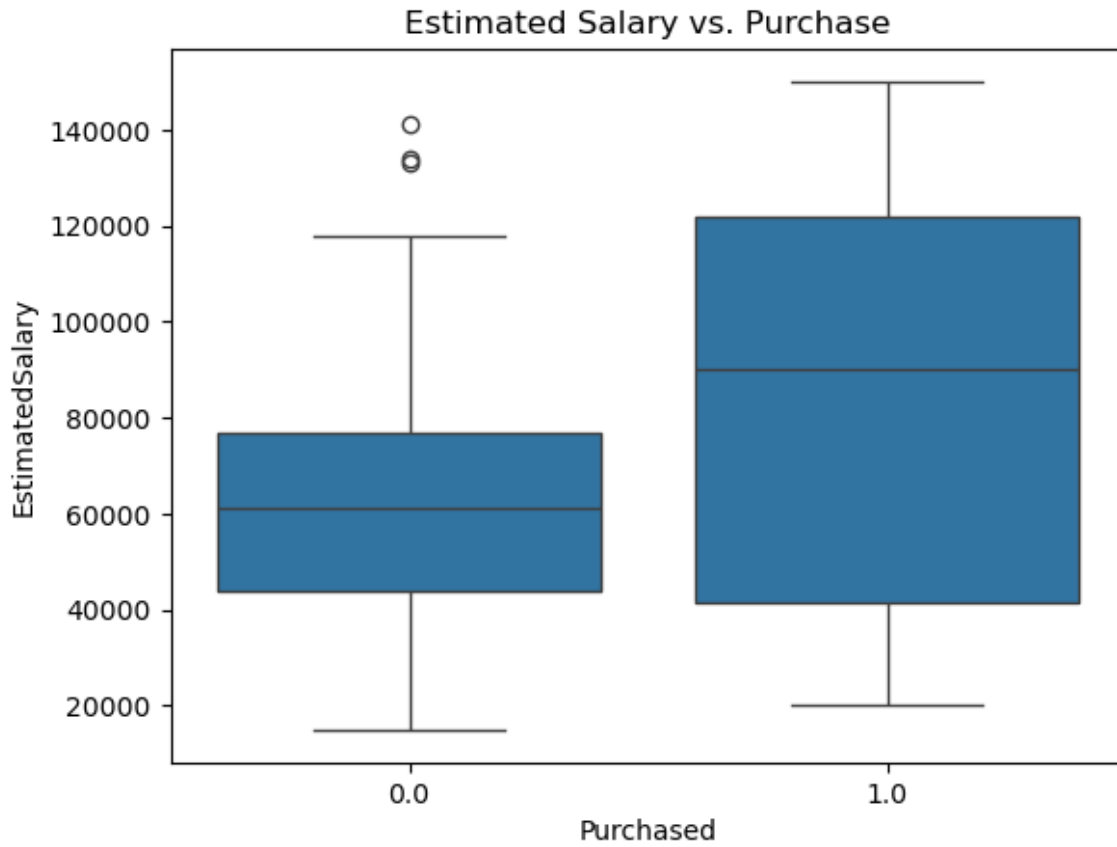
In [9]:

```
1 # Exploratory Questions:
2
3 # Is there a correlation between age and the likelihood of making a purchase?
4 sns.boxplot(x='Purchased', y='Age', data=data)
5 plt.title('Age vs. Purchase')
6 plt.show()
7
8
```



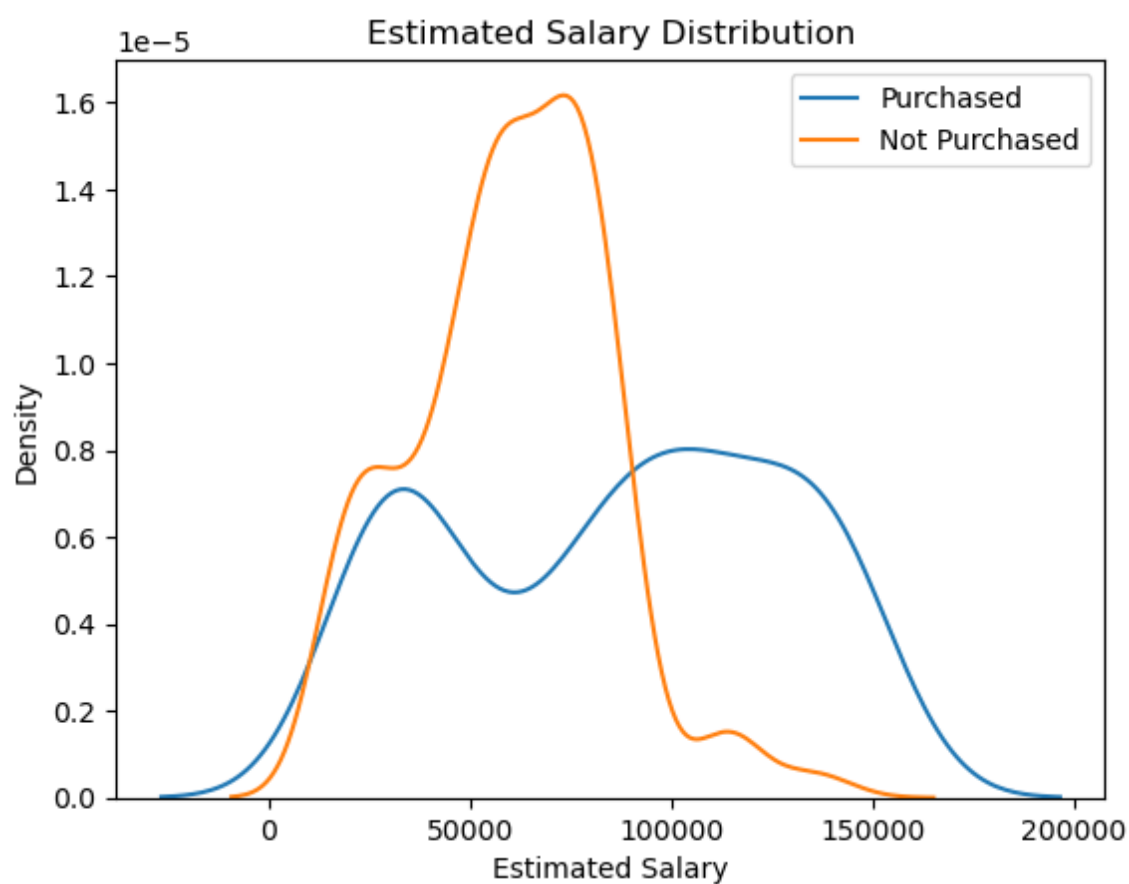
In [10]:

```
1 # Is there a correlation between estimated salary and the likelihood of making
2 sns.boxplot(x='Purchased', y='EstimatedSalary', data=data)
3 plt.title('Estimated Salary vs. Purchase')
4 plt.show()
5
6
```



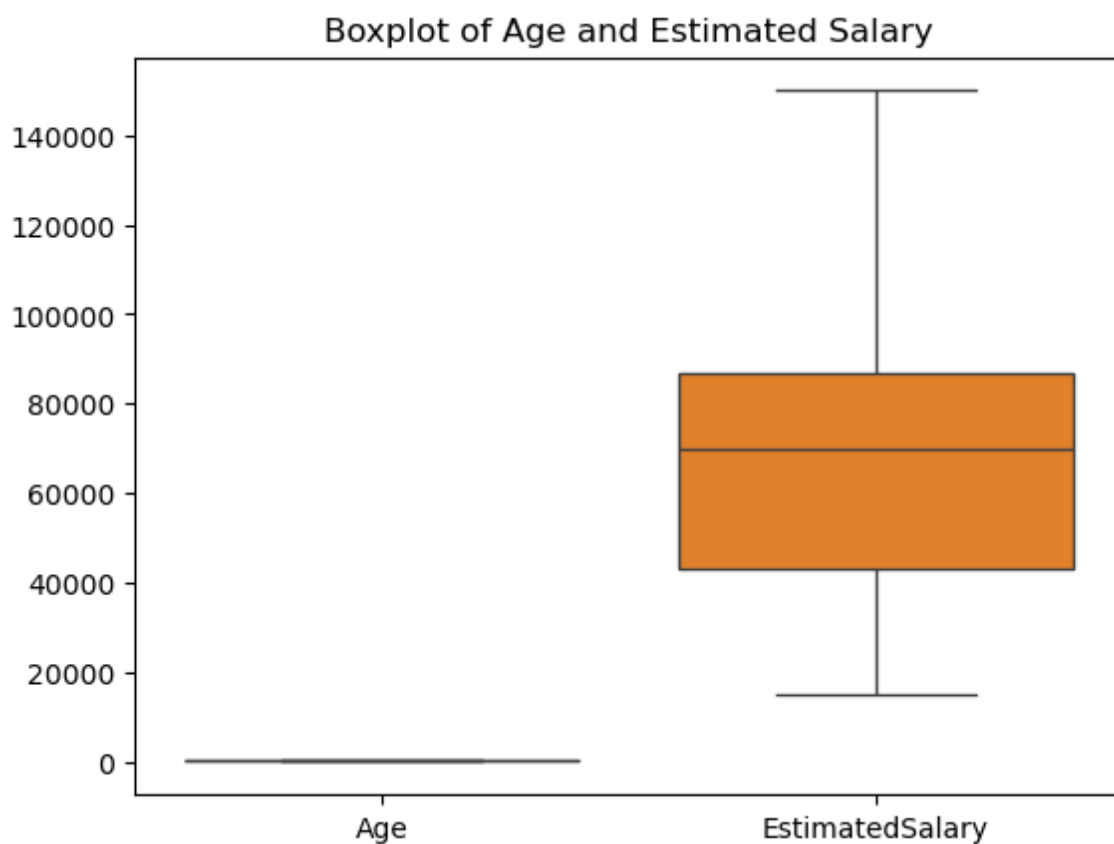
In [11]:

```
1 # How does the distribution of estimated salary differ between those who made
2 sns.kdeplot(data[data['Purchased'] == 1]['EstimatedSalary'], label='Purchased')
3 sns.kdeplot(data[data['Purchased'] == 0]['EstimatedSalary'], label='Not Purchased')
4 plt.title('Estimated Salary Distribution')
5 plt.xlabel('Estimated Salary')
6 plt.ylabel('Density')
7 plt.legend()
8 plt.show()
9
10
```



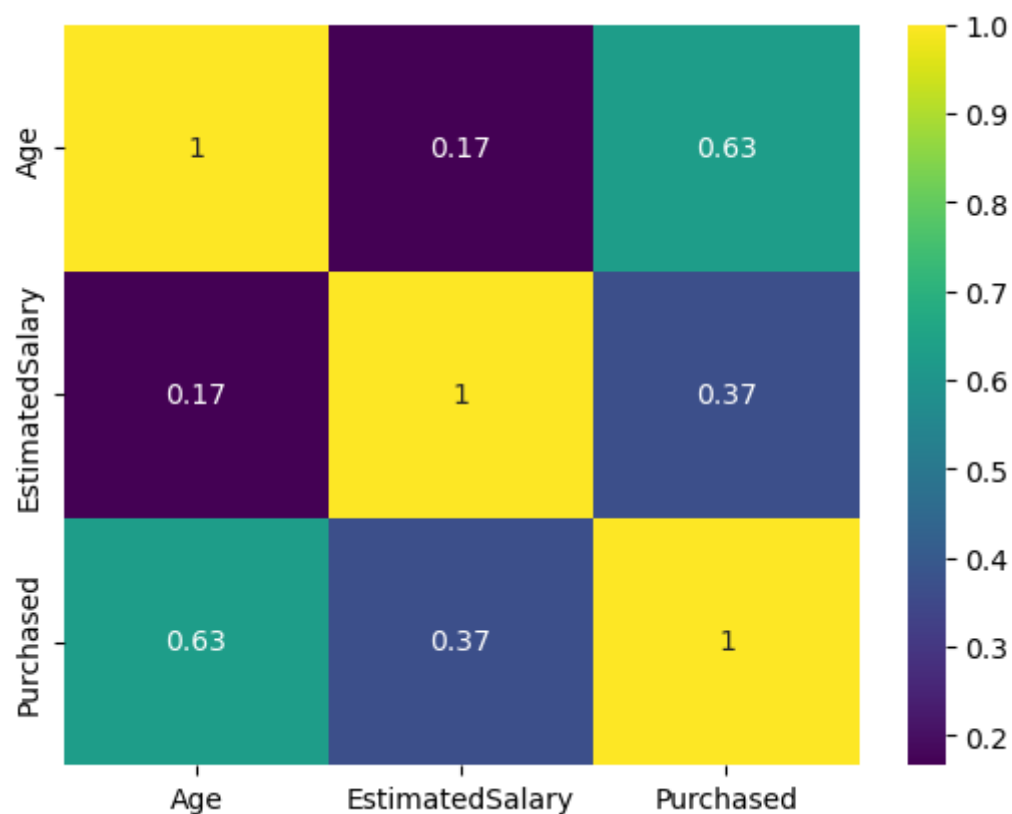
In [12]:

```
1 # Are there any outliers in the dataset, and how do they affect the analysis?
2 sns.boxplot(data=data[['Age', 'EstimatedSalary']])
3 plt.title('Boxplot of Age and Estimated Salary')
4 plt.show()
```



In [13]:

```
1 sns.heatmap(data.corr(), annot = True, cmap='viridis');
```



- the Age appears to have more impact in the purchased column when compared to Estimated salary with 63% positive relationship/correlations

In []:

```
1
```