

```
In [4]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.metrics import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.over_sampling import SMOTE
from sklearn.metrics.classifier import

In [5]: df = pd.read_csv('C:\Users\Adrian\OneDrive\Desktop\CustomerTravel.csv')
print(df)

Age FrequentFlyer AnnualIncomeClass ServicesOpted \
0 34 No Middle Income 6
1 34 Yes Low Income 5
2 37 No Middle Income 3
3 39 No Middle Income 2
4 39 No Low Income 1
... ..
949 31 Yes Low Income 1
950 39 No Middle Income 5
951 37 No Middle Income 4
952 39 No Low Income 1
953 31 Yes High Income 1

AccountSyncedToSocialMedia BookedHotelOrNot Target
0 0 No Yes 0
1 0 Yes No 1
2 0 Yes No 0
3 0 No No 0
4 0 No No 0
... ..
949 0 No No 0
950 0 No Yes 0
951 0 No No 0
952 0 Yes Yes 0
953 0 No No 0

[954 rows x 7 columns]

In [10]: print(df.isnull().sum())

# Steps 1-5: Handle missing values for numeric columns only
df.fillna(df.mean(numeric_only=True), inplace=True)

# Check for missing values after filling
print(df.isnull().sum())

Age FrequentFlyer AnnualIncomeClass ServicesOpted \
0 0
1 0
2 0
3 0
4 0
... ..
949 0
950 0
951 0
952 0
953 0

AccountSyncedToSocialMedia BookedHotelOrNot Target
0 0
1 0
2 0
3 0
4 0
... ..
949 0
950 0
951 0
952 0
953 0

dtype: int64

In [11]: label_encoder = LabelEncoder()
categorical_cols = ['FrequentFlyer', 'AnnualIncomeClass', 'AccountSyncedToSocialMedia', 'BookedHotelOrNot']
for col in categorical_cols:
    df[col] = label_encoder.fit_transform(df[col])

In [14]: scaler = StandardScaler()
df[['Age', 'ServicesOpted']] = scaler.fit_transform(df[['Age', 'ServicesOpted']])

# Display the preprocessed dataset
print(df.head())

# Check the balance of the target variable
sns.countplot(x='Target', data=df)
plt.title('Class Balance')
plt.show()

# Display the count of each class
print(df['Target'].value_counts())

# Split the data into features and target
X = df.drop('Target', axis=1)
y = df['Target']

Age FrequentFlyer AnnualIncomeClass ServicesOpted \
0 0.569884 0 0 2.219338
1 0.569884 2 1 1.596429
2 1.469282 0 2 0.396227
3 -0.632267 0 2 -0.272274
4 -0.632267 0 1 -0.895175

AccountSyncedToSocialMedia BookedHotelOrNot Target
0 0 1 0
1 0 1 0
2 1 0 1
3 0 0 0
4 0 0 0

Class Balance
count
Target
0 739
1 224
Name: count, dtype: int64

Import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.metrics.classifier import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.over_sampling import SMOTE
```

Load the new dataset

```
df = pd.read_csv('C:\Users\Adrian\OneDrive\Desktop\CustomerTravel.csv')
```

Handle missing values for numeric columns only

```
df.fillna(df.mean(numeric_only=True), inplace=True)
```

Encode categorical variables

```
label_encoder = LabelEncoder()
categorical_cols = ['FrequentFlyer', 'AnnualIncomeClass', 'AccountSyncedToSocialMedia', 'BookedHotelOrNot']
for col in categorical_cols:
    df[col] = label_encoder.fit_transform(df[col])
```

Feature scaling

```
scaler = StandardScaler()
df[['Age', 'ServicesOpted']] = scaler.fit_transform(df[['Age', 'ServicesOpted']])
```

Split the data into features and target

```
X = df.drop('Target', axis=1)
y = df['Target']
```

Split the data into training and testing sets

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Apply SMOTE to balance the classes in the training set

```
smote = SMOTE(random_state=42)
X_train_res, y_train_res = smote.fit_resample(X_train, y_train)
```

Plot class balance before and after SMOTE

```
fig, ax = plt.subplots(1, 2, figsize=(12, 6))
```

Before SMOTE

```
sns.countplot(X_train_res, ax=ax[0], x=df['set_label('Class Balance Before SMOTE')'], ax[0].set_ylabel('Class') and[0].set_ylabel('Count')
```

After SMOTE

```
sns.countplot(X_train_res, ax=ax[1], x=df['set_label('Class Balance After SMOTE')'], ax[1].set_ylabel('Class') and[1].set_ylabel('Count')
```

```
plt.tight_layout()
plt.show()
```

```
In [18]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Apply SMOTE to balance the classes in the training set
smote = SMOTE(random_state=42)
X_train_res, y_train_res = smote.fit_resample(X_train, y_train)

# Check the new balance
sns.countplot(X_train_res, ax=ax[0], x=df['set_label('Class Balance Before SMOTE')'], ax[0].set_ylabel('Class') and[0].set_ylabel('Count')
```

```
In [24]: print(df.columns)

Index(['Age', 'FrequentFlyer', 'AnnualIncomeClass', 'ServicesOpted',
       'AccountSyncedToSocialMedia', 'BookedHotelOrNot', 'Target'],
      dtype='object')

In [26]: from sklearn.decomposition import PCA
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Print column names to check if 'Cluster' exists
print('Columns in DataFrame:', df.columns)

# Check if 'Cluster' is in the DataFrame
if 'Cluster' in df.columns:
    pca = PCA(n_components=2)
    df_pca = pca.fit_transform(df.drop('Cluster', axis=1))
    df_pca = pd.DataFrame(df_pca, columns=['PC1', 'PC2'])
    df_pca['Cluster'] = df['Cluster']

# Visualize the clusters
plt.figure(figsize=(10, 8))
sns.scatterplot(x=df_pca['PC1'], y=df_pca['PC2'], hue='Cluster', data=df_pca, palette='viridis', s=300)
plt.title('Customer Segmentation using K-means Clustering')
plt.xlabel('PC1')
plt.ylabel('PC2')
plt.legend('Cluster')
plt.grid(True)
plt.savefig('k-means_clusters.png', dpi=300, bbox_inches='tight')
plt.show()

else:
    print("Cluster column not found in DataFrame.")

columns in DataFrame: Index(['Age', 'FrequentFlyer', 'AnnualIncomeClass', 'ServicesOpted',
       'AccountSyncedToSocialMedia', 'BookedHotelOrNot', 'Target'],
      dtype='object')
Cluster column not found in DataFrame.
```

```
In [28]: print(df.head())

Age FrequentFlyer AnnualIncomeClass ServicesOpted \
0 0.569884 0 2 2.219338
1 0.569884 2 1 1.596429
2 1.469282 0 2 0.396227
3 -0.632267 0 2 -0.272274
4 -0.632267 0 1 -0.895175

AccountSyncedToSocialMedia BookedHotelOrNot Target
0 0 1 0
1 0 1 0
2 1 0 1
3 0 0 0
4 0 0 0
```

```
In [43]: import pandas as pd

# Example: Loading a DataFrame from a CSV file
churn_customers = pd.read_csv('C:\Users\Adrian\OneDrive\Desktop\CustomerTravel.csv')

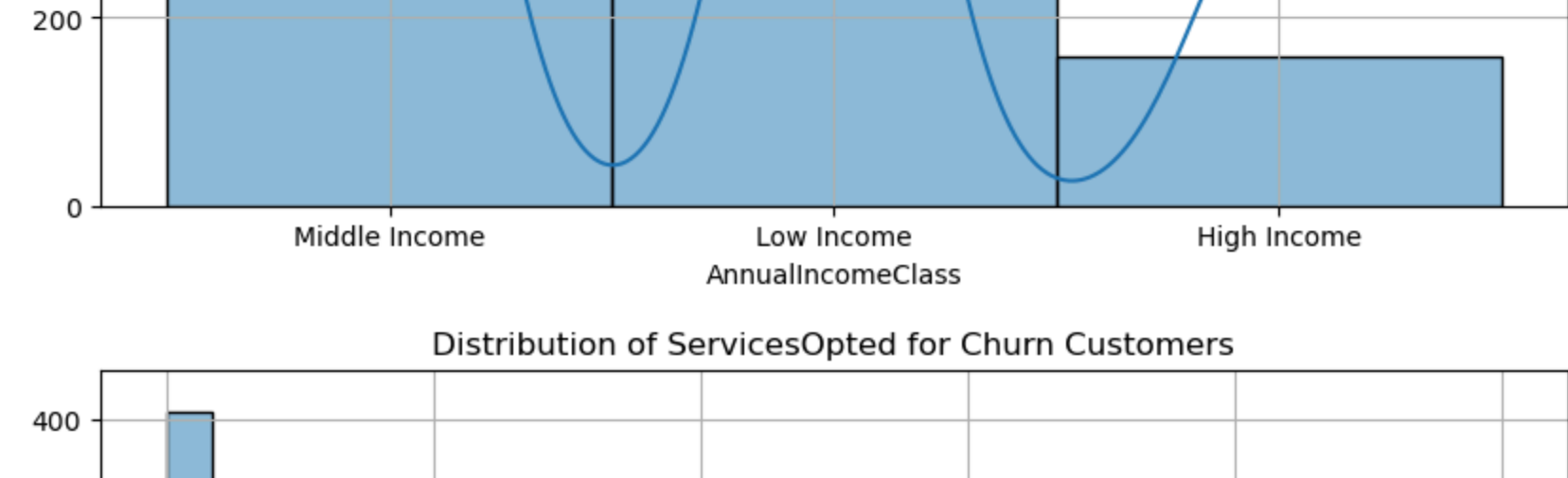
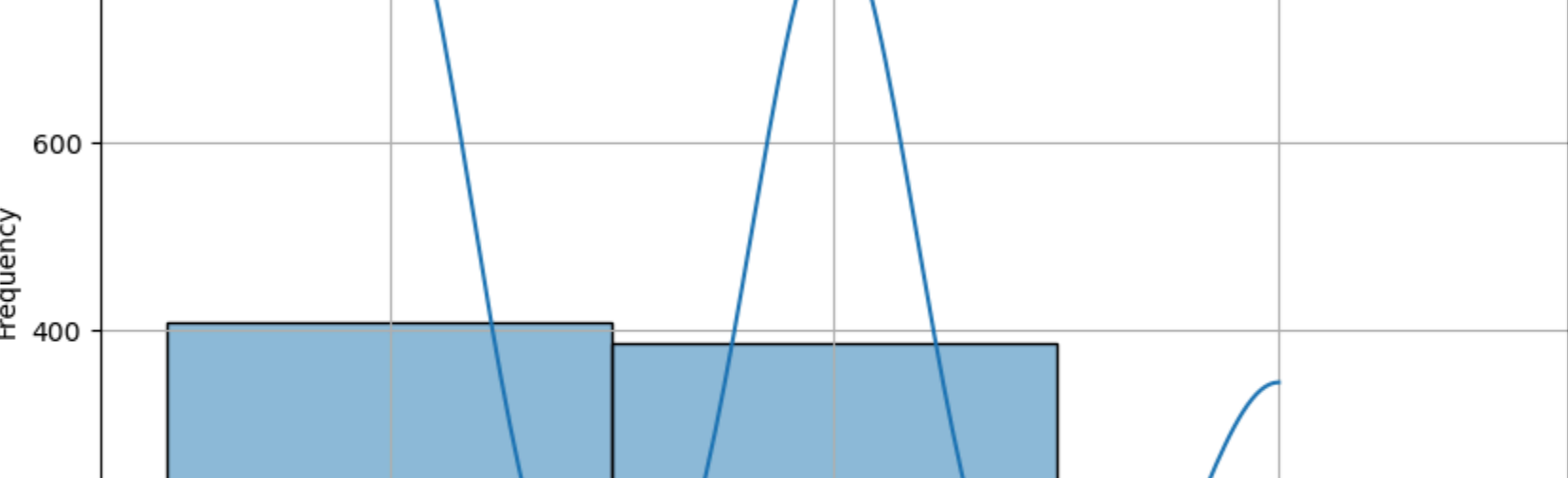
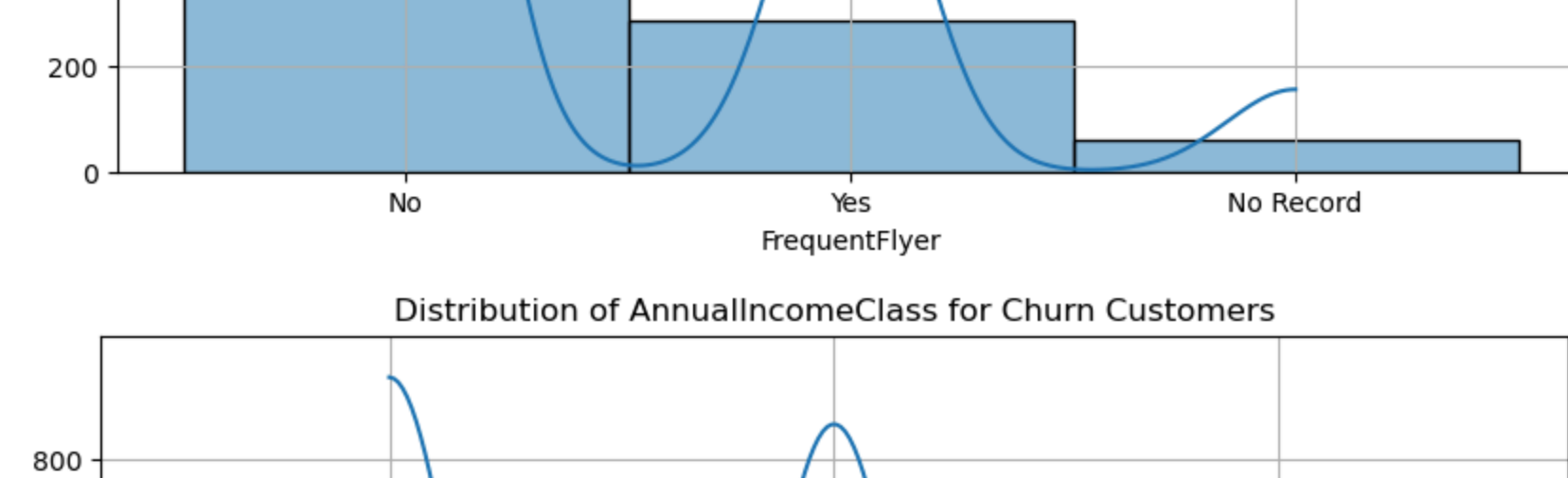
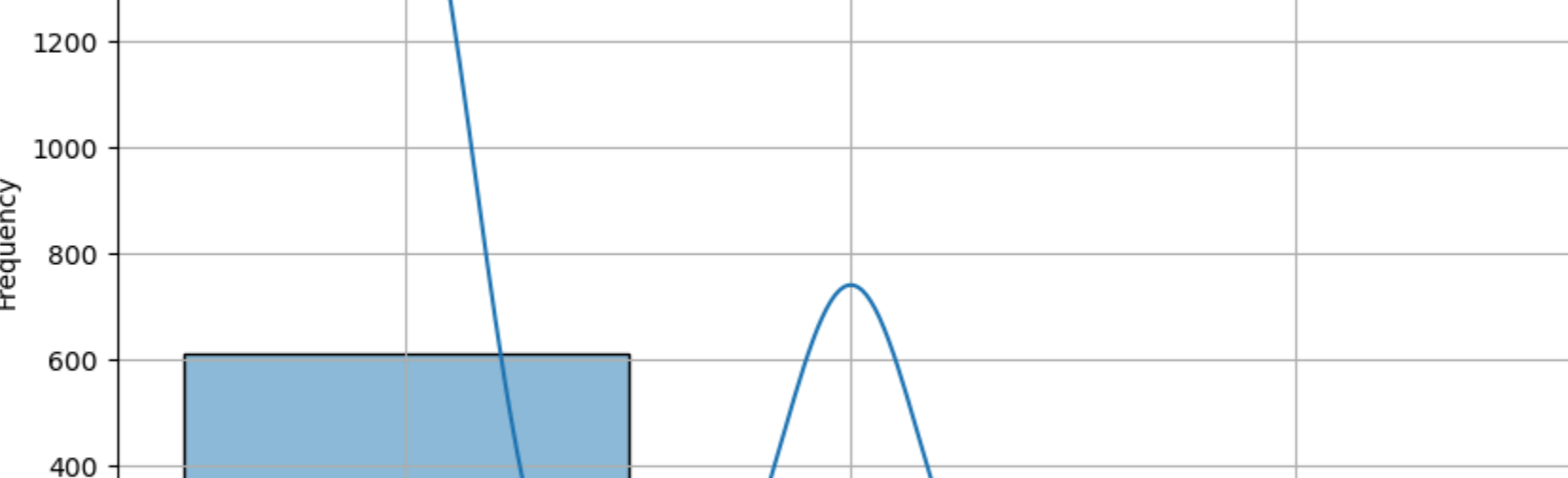
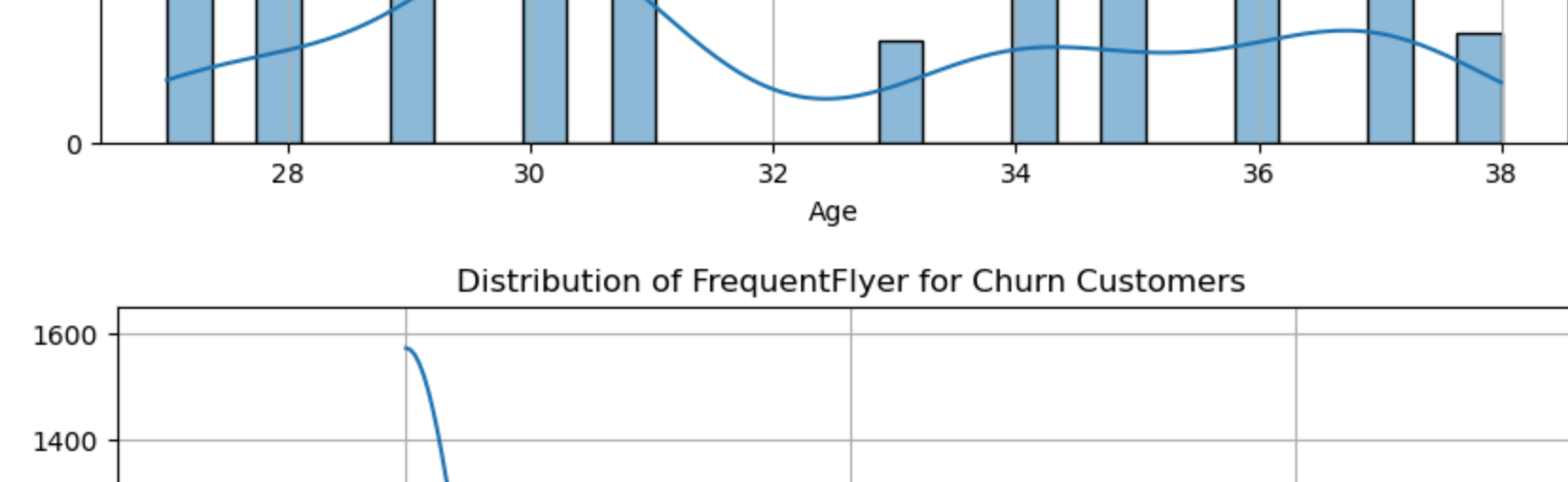
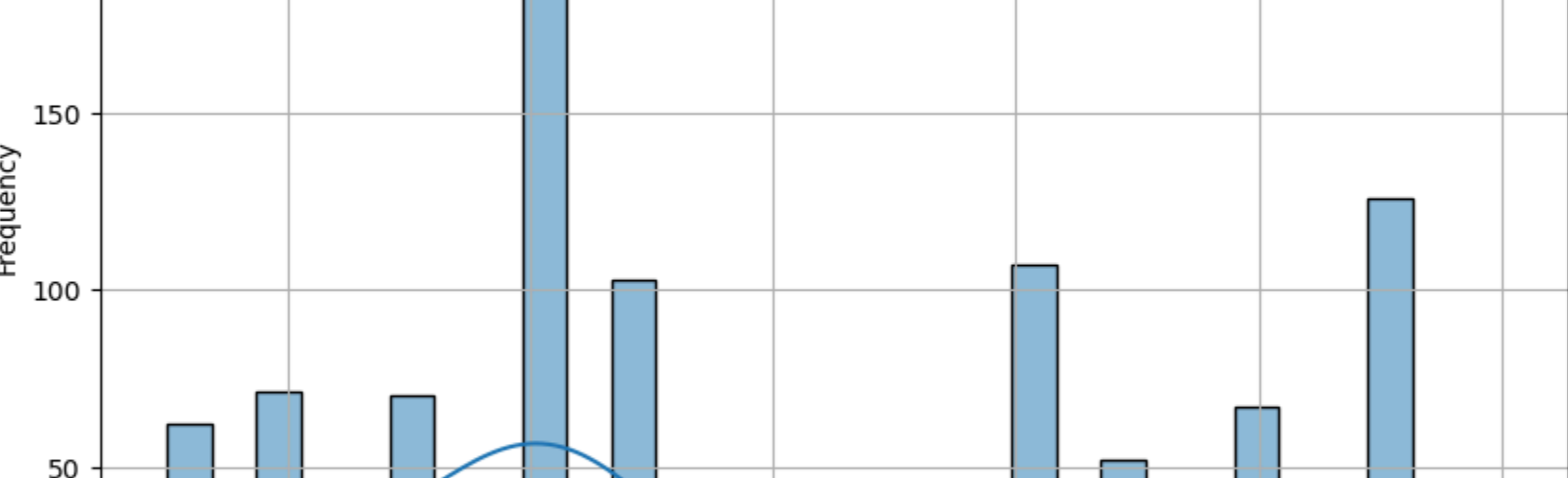
# Verify DataFrame is loaded
print(churn_customers.head())

# List of important features
important_features = [
    'Age',
    'FrequentFlyer',
    'AnnualIncomeClass',
    'ServicesOpted',
    'AccountSyncedToSocialMedia',
    'BookedHotelOrNot'
]

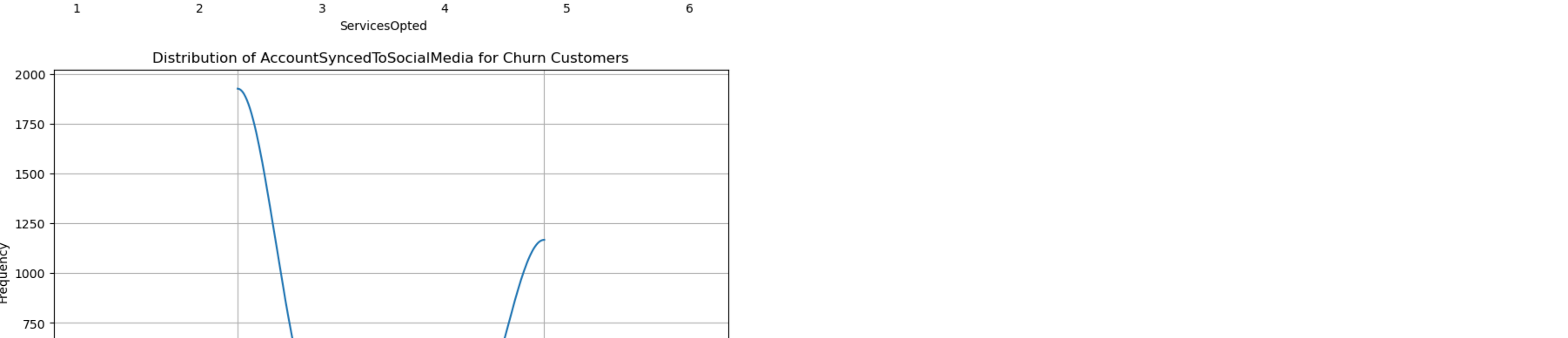
# Visualize important features
for feature in important_features:
    if feature in churn_customers.columns:
        sns.histplot(churn_customers[feature], kde=True, bins=30)
        plt.xlabel(feature)
        plt.ylabel('Frequency')
        plt.grid(True)
        plt.show()
    else:
        print(f"Feature '{feature}' not found in DataFrame.")

Age FrequentFlyer AnnualIncomeClass ServicesOpted \
0 0.569884 0 2 2.219338
1 0.569884 2 1 1.596429
2 1.469282 0 2 0.396227
3 -0.632267 0 2 -0.272274
4 -0.632267 0 1 -0.895175

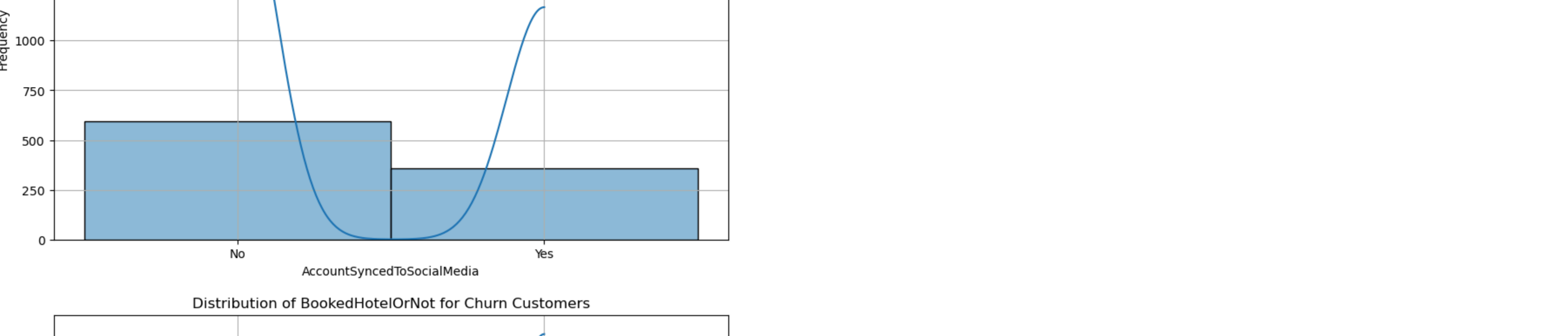
AccountSyncedToSocialMedia BookedHotelOrNot Target
0 0 1 0
1 0 1 0
2 1 0 1
3 0 0 0
4 0 0 0
```



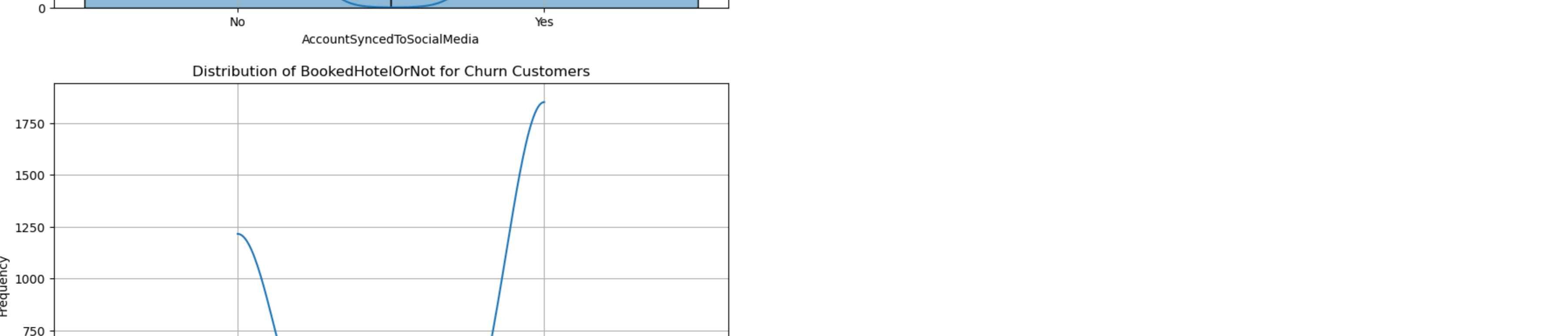
```
In [46]: import matplotlib.pyplot as plt
import numpy as np
fig = plt.figure(figsize=(10, 8))
df.groupby('AnnualIncomeClass').Target.value_counts().unstack().plot(kind='bar', ax=fig)
plt.show()
```



```
In [48]: fig = plt.figure(figsize=(10, 8))
df.groupby('AnnualIncomeClass').FrequentFlyer.value_counts().unstack().plot(kind='bar', ax=fig)
plt.show()
```



```
In [50]: fig = plt.figure(figsize=(10, 8))
df.groupby('AnnualIncomeClass').ServicesOpted.value_counts().unstack().plot(kind='bar', ax=fig)
plt.show()
```

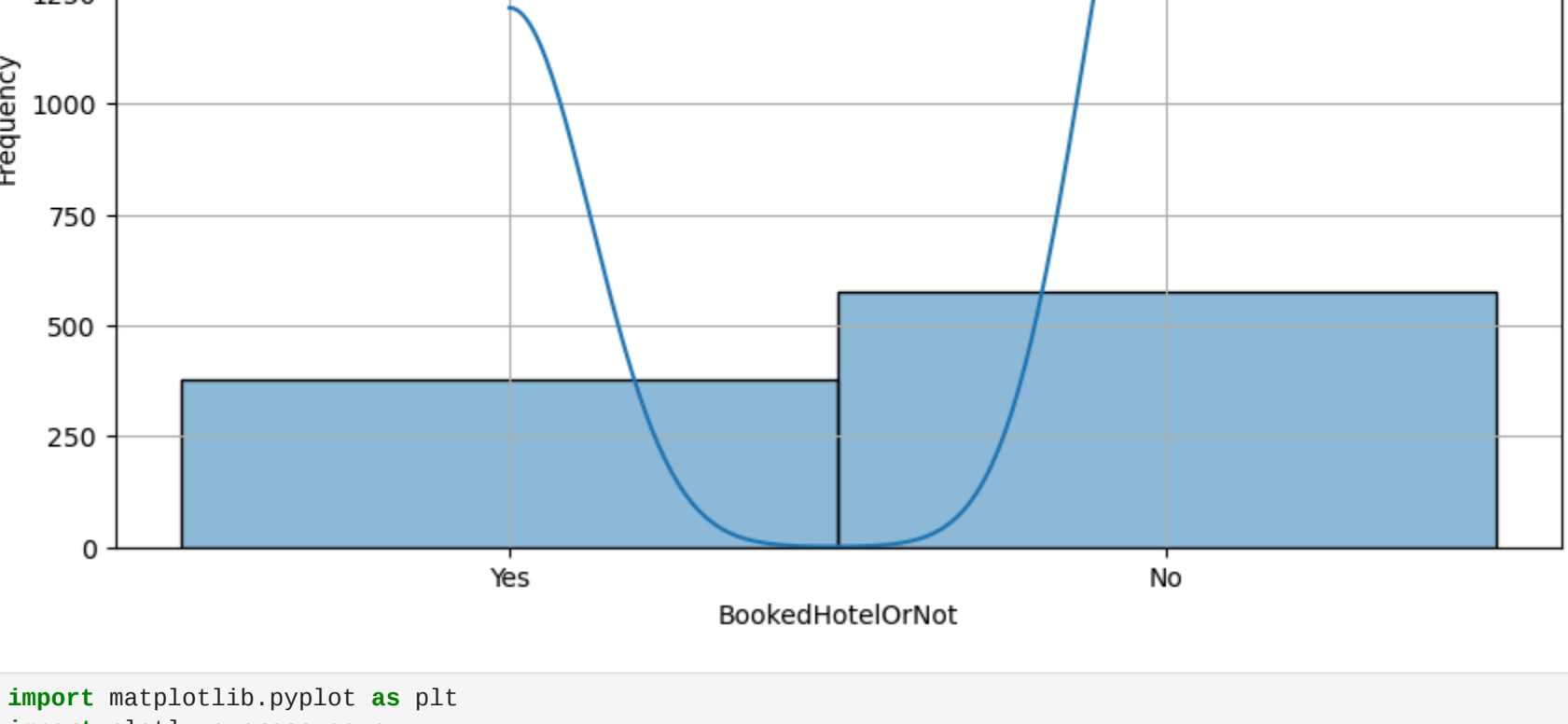


```
In [54]: import matplotlib.pyplot as plt

# Example data for the pie chart
income_data = pd.Series([0, 45, 25], index=['High Income', 'Low Income', 'Middle Income'])

# Create the pie chart
plt.pie(income_data, labels=income_data.index, autopct='%1.1f%%', radius=1.2, textprops={'fontsize': 16})

plt.title('How Income Impacts Customer Churn', color='b')
plt.show()
```



```
In [63]: churn_perc = df['BookedHotelOrNot'].value_counts(normalize=1)
print(churn_perc)

plt.figure(figsize=(5, 6))
sns.pieplot(churn_perc, df['BookedHotelOrNot'].value_counts(normalize=1), edgecolor='black', width=0.8, color=['skyblue', 'purple'])
plt.title('How Hotel Bookings Affect Customer Churn')
plt.show()
```



```
In [67]: print('Statistical Summary')
df.describe()

Out[67]:
```

	count	mean	std	min	25%	50%	75%	max
Age	954.0	0.468317e+16	1.000325	-1.531645	-0.632267	-0.332475	0.866696	1.766075
FrequentFlyer	954.0	6.624738e-01	0.907770	0.000000	0.000000	0.000000	2.000000	2.000000
AnnualIncomeClass	954.0	1.262050e+00	0.728132	0.000000	1.000000	1.000000	2.000000	2.000000
ServicesOpted	954.0	3.351917e-17	1.000325	-0.895175	-0.895175	-0.272274	0.973538	2.219338
AccountSyncedToSocialMedia	954.0	3.779856e-01	0.484969	0.000000	0.000000	0.000000	1.000000	1.000000
BookedHotelOrNot	954.0	1.962654e-01	0.489369	0.000000	0.000000	0.000000	1.000000	1.000000
Target	954.0	2.348000e-01	0.424597	0.000000	0.000000	0.000000	0.000000	1.000000

```
In [72]: print(df)

Age FrequentFlyer AnnualIncomeClass ServicesOpted \
0 0.569884 0 0 2.219338
1 0.569884 2 1 1.596429
2 1.469282 0 2 0.396227
3 -0.632267 0 2 -0.272274
4 -0.632267 0 1 -0.895175
... ..
949 0.332475 2 1 -0.895175
950 -0.632267 0 2 1.596429
951 1.469282 0 2 0.396227
952 -0.632267 0 1 -0.895175
953 -0.332475 2 0 0.000000

AccountSyncedToSocialMedia BookedHotelOrNot Target
0 0 1 0
1 0 1 0
2 1 0 1
3 0 0 0
4 0 0 0
... ..
949 0 0 0
950 0 1 0
951 0 0 0
952 1 1 0
953 0 0 0

[954 rows x 7 columns]
```

```
In [74]: print(churn_perc)

BookedHotelOrNot
0 0.603714
1 0.396286
Name: proportion, dtype: float64
```

