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

```
pip install pymysql sqlalchemy pandas
```

Requirement already satisfied: pymysql in c:\users\admin\appdata\local\progr
ams\python\python313\lib\site-packages (1.1.1)
Requirement already satisfied: sqlalchemy in c:\users\admin\appdata\local\pr
ograms\python\python313\lib\site-packages (2.0.41)
Requirement already satisfied: pandas in c:\users\admin\appdata\local\progra
ms\python\python313\lib\site-packages (2.3.0)
Requirement already satisfied: greenlet>=1 in c:\users\admin\appdata\local\p
rograms\python\python313\lib\site-packages (from sqlalchemy) (3.2.3)
Requirement already satisfied: typing-extensions>=4.6.0 in c:\users\admin\ap
pdata\local\programs\python\python313\lib\site-packages (from sqlalchemy)
(4.14.0)
Requirement already satisfied: numpy>=1.26.0 in c:\users\admin\appdata\local
\programs\python\python313\lib\site-packages (from pandas) (2.3.1)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\users\admin\appd
ata\local\programs\python\python313\lib\site-packages (from pandas) (2.9.0.p
ost0)
Requirement already satisfied: pytz>=2020.1 in c:\users\admin\appdata\local
\programs\python\python313\lib\site-packages (from pandas) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in c:\users\admin\appdata\loca
l\programs\python\python313\lib\site-packages (from pandas) (2025.2)
Requirement already satisfied: six>=1.5 in c:\users\admin\appdata\local\prog
rams\python\python313\lib\site-packages (from python-dateutil>=2.8.2->panda
s) (1.17.0)
Note: you may need to restart the kernel to use updated packages.

```python
from sqlalchemy import create_engine
import pandas as pd

user = "root"
password = "Aryan%402004"  # '@' becomes '%40'
host = "localhost"
port = 3306
database = "banking_case"  # replace with your actual DB name

engine = create_engine(f"mysql+pymysql://{user}:{password}@{host}:{port}/{da

# Test query
df = pd.read_sql("SHOW TABLES", engine)
print(df)
```

```
  Tables_in_banking_case
0              customer
```

```python
query= "SELECT * FROM banking_case.customer"
```

```python
df= pd.read_sql(query,engine)
df.head(5)
```

Out[6]:

| | ï»¿Client ID | Name | Age | Location ID | Joined Bank | Banking Contact | Nationality | Occupation |
|---|---|---|---|---|---|---|---|---|
| **0** | IND81288 | Raymond Mills | 24 | 34324 | 06-05-2019 | Anthony Torres | American | Safety Technician I' |
| **1** | IND65833 | Julia Spencer | 23 | 42205 | 10-12-2001 | Jonathan Hawkins | African | Software Consultan |
| **2** | IND47499 | Stephen Murray | 27 | 7314 | 25-01-2010 | Anthony Berry | European | Help Des Operato |
| **3** | IND72498 | Virginia Garza | 40 | 34594 | 28-03-2019 | Steve Diaz | American | Geologist |
| **4** | IND60181 | Melissa Sanders | 46 | 41269 | 20-07-2012 | Shawn Long | American | Assistan Professo |

5 rows × 25 columns

In [7]:
```python
#Generate Descriptive statistics for the database
df.describe()
```

Out[7]:

| | Age | Location ID | Estimated Income | Superannuation Savings | Amount of Credit Cards |
|---|---|---|---|---|---|
| **count** | 3000.000000 | 3000.000000 | 3000.000000 | 3000.000000 | 3000.000000 |
| **mean** | 51.039667 | 21563.323000 | 171305.034263 | 25531.599673 | 1.463667 |
| **std** | 19.854760 | 12462.273017 | 111935.808209 | 16259.950770 | 0.676387 |
| **min** | 17.000000 | 12.000000 | 15919.480000 | 1482.030000 | 1.000000 |
| **25%** | 34.000000 | 10803.500000 | 82906.595000 | 12513.775000 | 1.000000 |
| **50%** | 51.000000 | 21129.500000 | 142313.480000 | 22357.355000 | 1.000000 |
| **75%** | 69.000000 | 32054.500000 | 242290.305000 | 35464.740000 | 2.000000 |
| **max** | 85.000000 | 43369.000000 | 522330.260000 | 75963.900000 | 3.000000 |

In [8]:
```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3000 entries, 0 to 2999
Data columns (total 25 columns):
 #   Column                   Non-Null Count  Dtype
---  ------                   --------------  -----
 0   ï»¿Client ID             3000 non-null   object
 1   Name                     3000 non-null   object
 2   Age                      3000 non-null   int64
 3   Location ID              3000 non-null   int64
 4   Joined Bank              3000 non-null   object
 5   Banking Contact          3000 non-null   object
 6   Nationality              3000 non-null   object
 7   Occupation               3000 non-null   object
 8   Fee Structure            3000 non-null   object
 9   Loyalty Classification   3000 non-null   object
 10  Estimated Income         3000 non-null   float64
 11  Superannuation Savings   3000 non-null   float64
 12  Amount of Credit Cards   3000 non-null   int64
 13  Credit Card Balance      3000 non-null   float64
 14  Bank Loans               3000 non-null   float64
 15  Bank Deposits            3000 non-null   float64
 16  Checking Accounts        3000 non-null   float64
 17  Saving Accounts          3000 non-null   float64
 18  Foreign Currency Account 3000 non-null   float64
 19  Business Lending         3000 non-null   float64
 20  Properties Owned         3000 non-null   int64
 21  Risk Weighting           3000 non-null   int64
 22  BRId                     3000 non-null   int64
 23  GenderId                 3000 non-null   int64
 24  IAId                     3000 non-null   int64
dtypes: float64(9), int64(8), object(8)
memory usage: 586.1+ KB
```

In [10]: `df.shape`

Out[10]:  (3000, 25)

In [12]:
```python
bins=[0,100000,300000, float('inf')]
labels=['low', 'med', 'high']

df['Income Band'] =pd.cut(df['Estimated Income'], bins=bins, labels=labels,
```
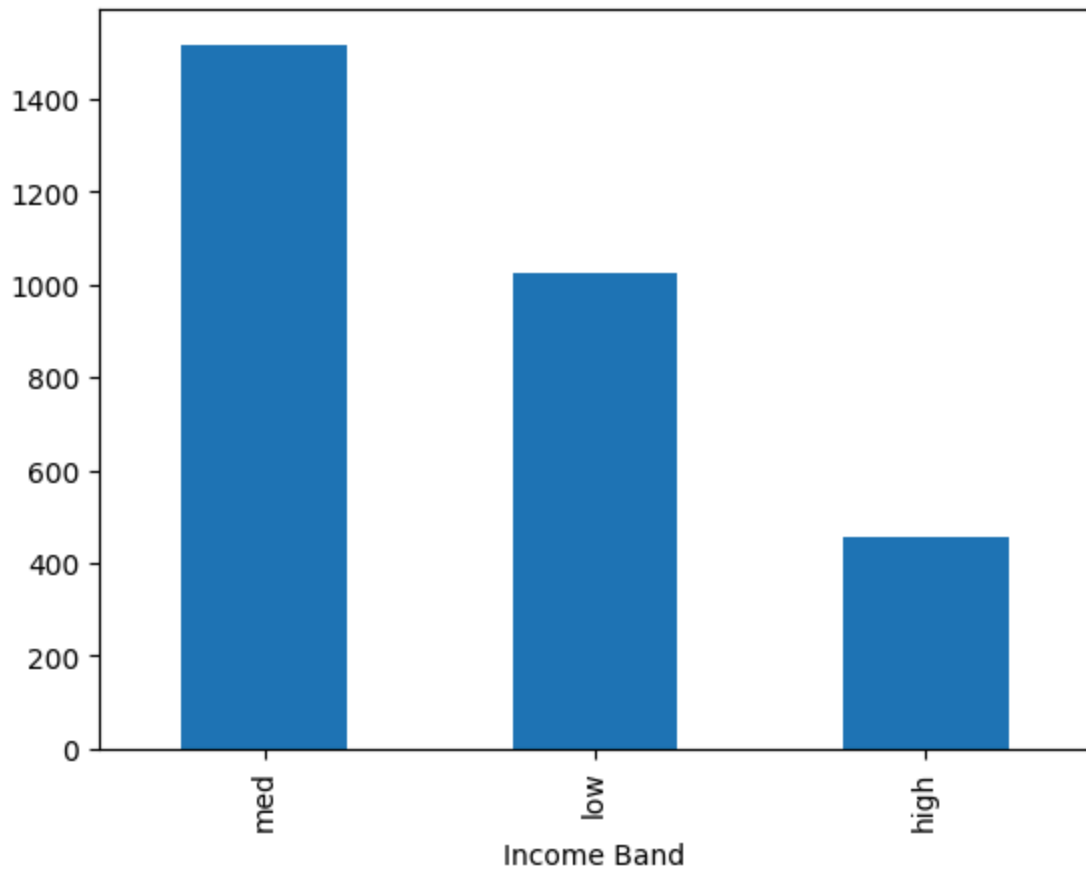
In [13]: `df['Income Band'].value_counts().plot(kind='bar')`

Out[13]:  <Axes: xlabel='Income Band'>

```
#Examine the distribution of unique categories in categorical columns
categorical_cols = df[["BRId", "GenderId","IAId","Amount of Credit Cards","M
for col in categorical_cols:
    print(f"Value Counts for '{col}':")
    display(df[col].value_counts())
```

```
Value Counts for 'BRId':
BRId
3    1352
1     660
2     495
4     493
Name: count, dtype: int64
Value Counts for 'GenderId':
GenderId
2    1512
1    1488
Name: count, dtype: int64
Value Counts for 'IAId':
```

```
IAId
1       177
2       177
3       177
4       177
8       177
9       176
13      176
12      176
10      176
11      176
14      176
15      176
6        89
5        89
7        89
16       88
17       88
18       88
19       88
20       88
21       88
22       88
Name: count, dtype: int64
Value Counts for 'Amount of Credit Cards':
Amount of Credit Cards
1    1922
2     765
3     313
Name: count, dtype: int64
Value Counts for 'Nationality':
Nationality
European      1309
Asian          754
American       507
Australian     254
African        176
Name: count, dtype: int64
Value Counts for 'Occupation':
Occupation
Associate Professor          28
Structural Analysis Engineer 28
Recruiter                    25
Account Coordinator          24
Human Resources Manager      24
                             ..
Office Assistant IV           8
Automation Specialist I       7
Computer Systems Analyst I    6
Developer III                 5
Senior Sales Associate        4
Name: count, Length: 195, dtype: int64
Value Counts for 'Fee Structure':
```

```
Fee Structure
High      1476
Mid        962
Low        562
Name: count, dtype: int64
Value Counts for 'Loyalty Classification':
Loyalty Classification
Jade         1331
Silver        767
Gold          585
Platinum      317
Name: count, dtype: int64
Value Counts for 'Properties Owned':
Properties Owned
2     777
1     776
3     742
0     705
Name: count, dtype: int64
Value Counts for 'Risk Weighting':
Risk Weighting
2     1222
1      836
3      460
4      322
5      160
Name: count, dtype: int64
Value Counts for 'Income Band':
Income Band
med      1517
low      1027
high      456
Name: count, dtype: int64
```

# Univariate Analysis

```python
In [29]:  sns.set(style="whitegrid")

          for col in categorical_cols:
              plt.figure(figsize=(6, 6))

              # Try to sort the categories numerically if possible
              try:
                  categories = sorted(df[col].dropna().unique(), key=lambda x: int(x))
              except:
                  categories = df[col].value_counts().index

              # Draw countplot
              ax = sns.countplot(data=df, x=col, order=categories, color="#4c72b0")

              # Add value labels on top of bars
              for p in ax.patches:
                  height = int(p.get_height())
                  ax.annotate(f'{height}', (p.get_x() + p.get_width() / 2., height),
```

```
                    ha='center', va='bottom', fontsize=10, color='black')

# Final formatting
plt.title(f"Count Plot for '{col}'", fontsize=14)
plt.xlabel(col, fontsize=12)
plt.ylabel("Count", fontsize=12)
plt.xticks(rotation=90)
plt.tight_layout()
plt.show()
```
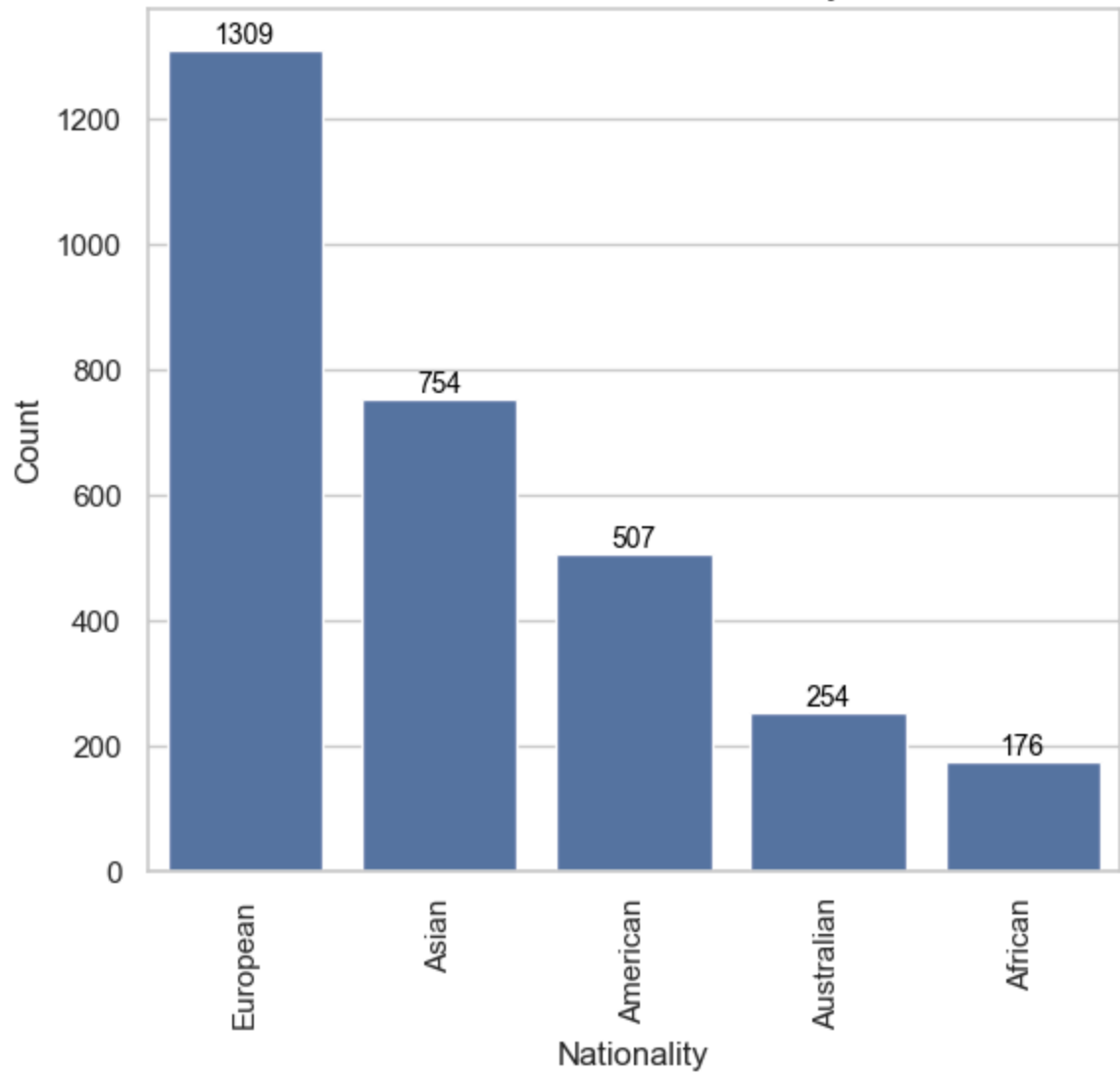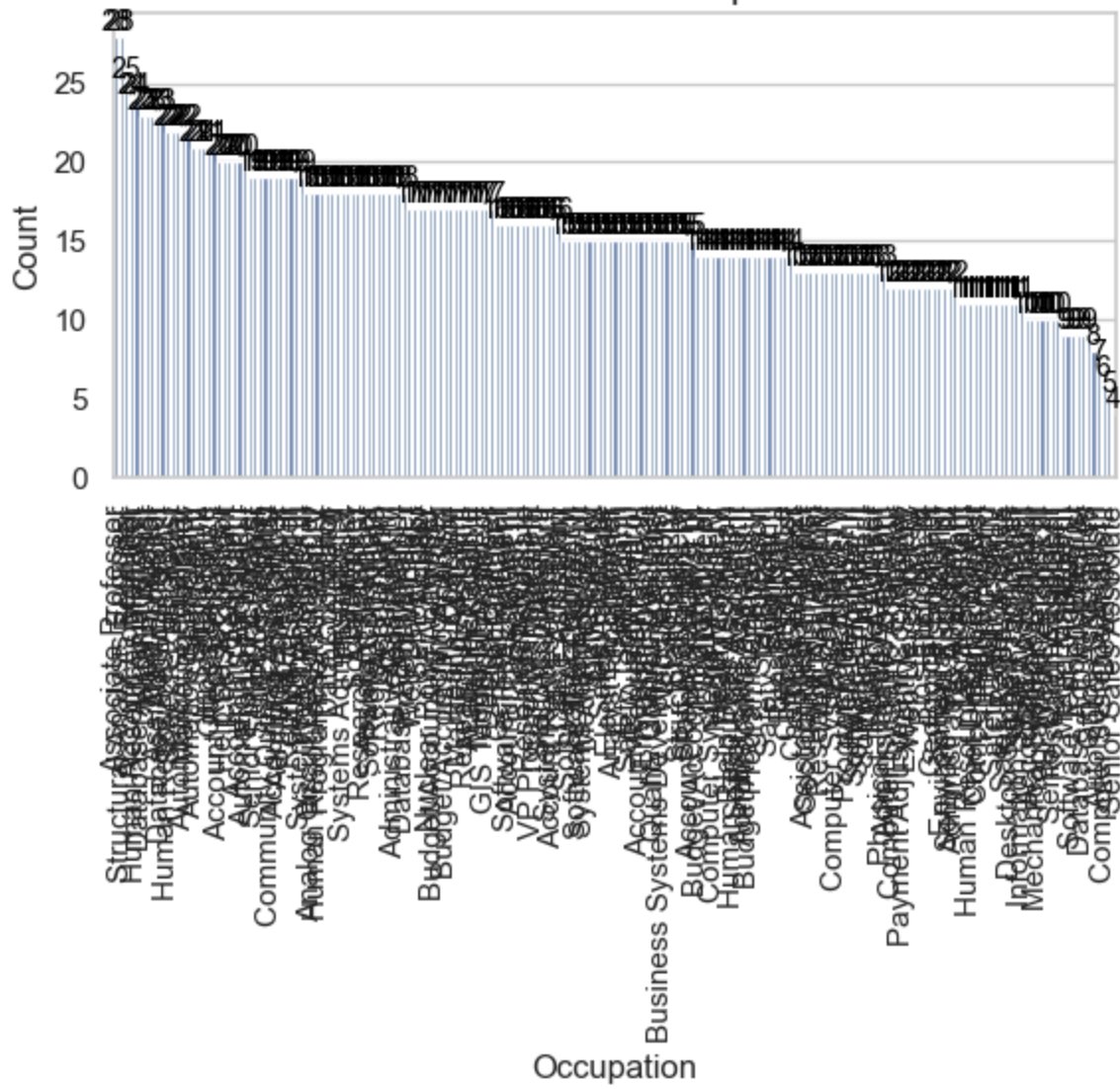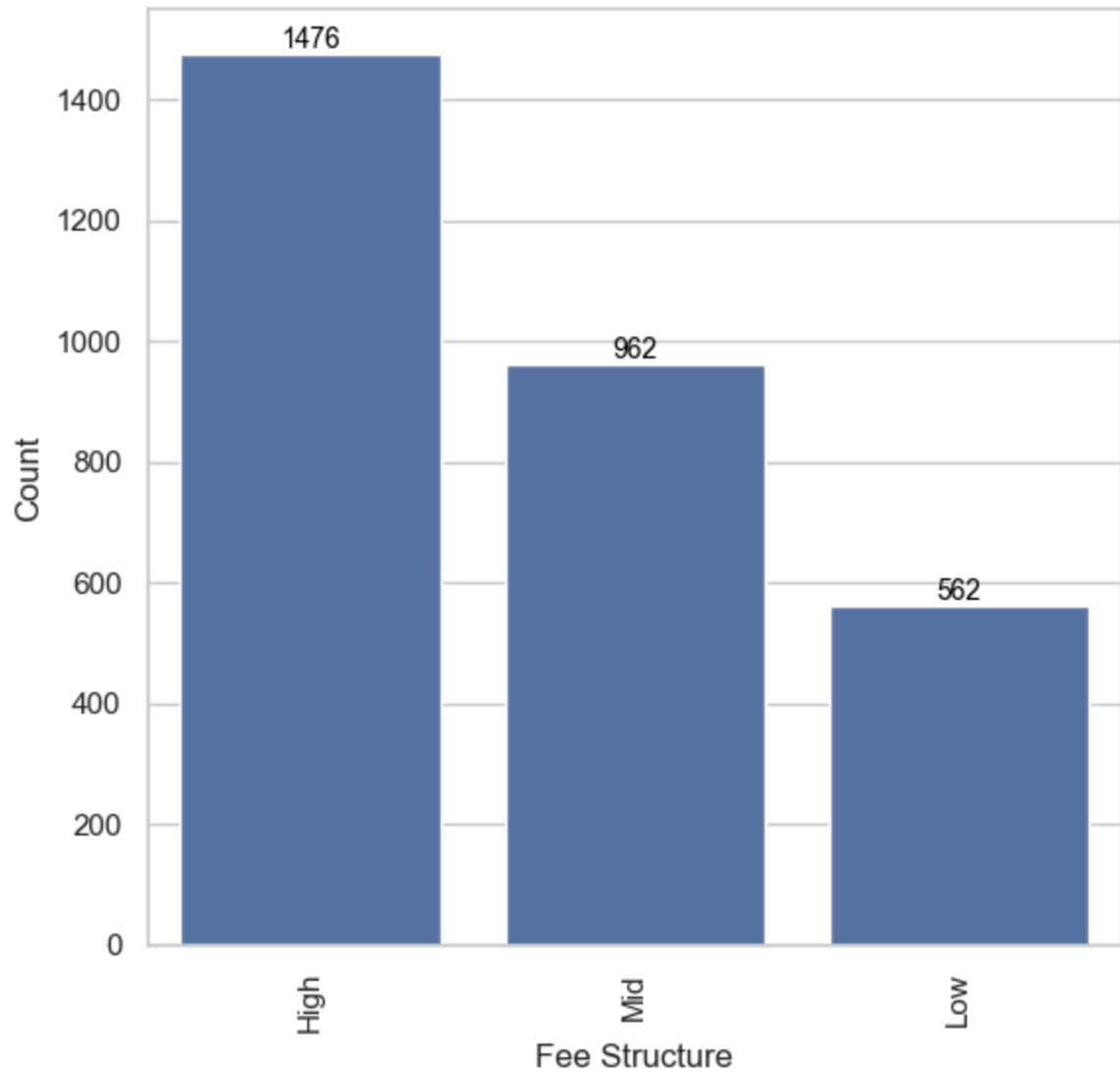


Count Plot for 'BRId'

Count Plot for 'GenderId'

Count Plot for 'lAld'

Count Plot for 'Amount of Credit Cards'
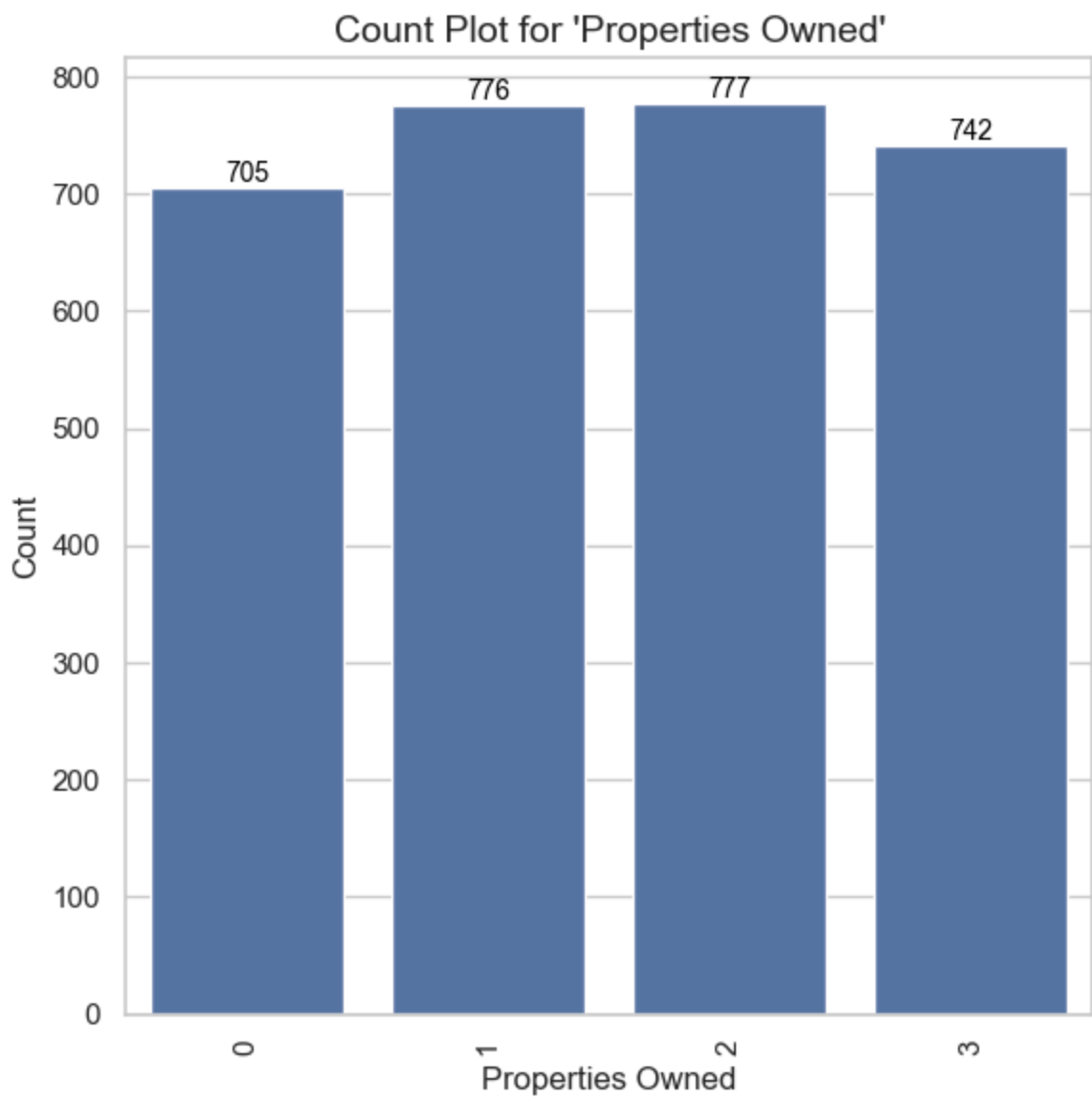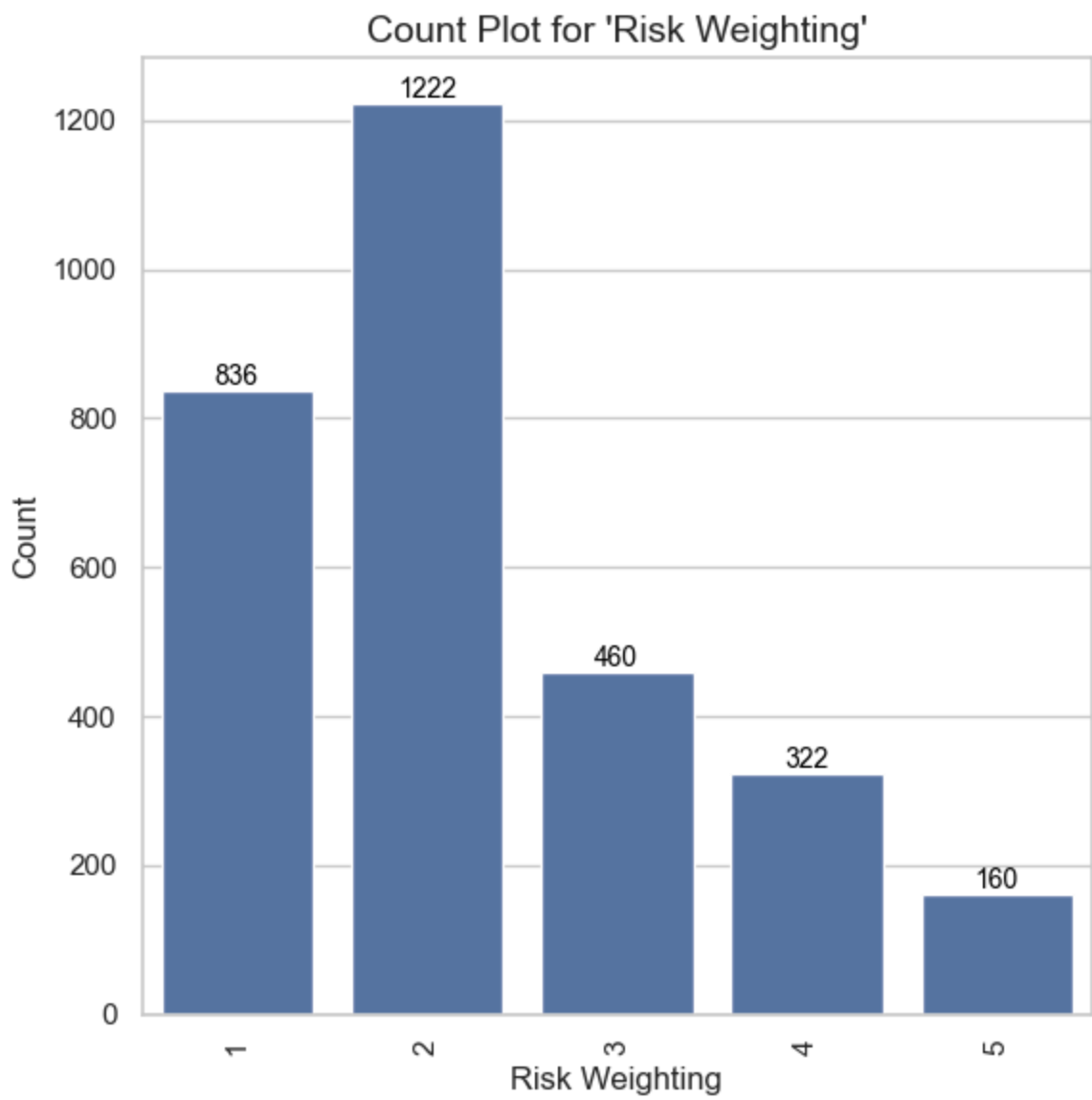
Count Plot for 'Nationality'
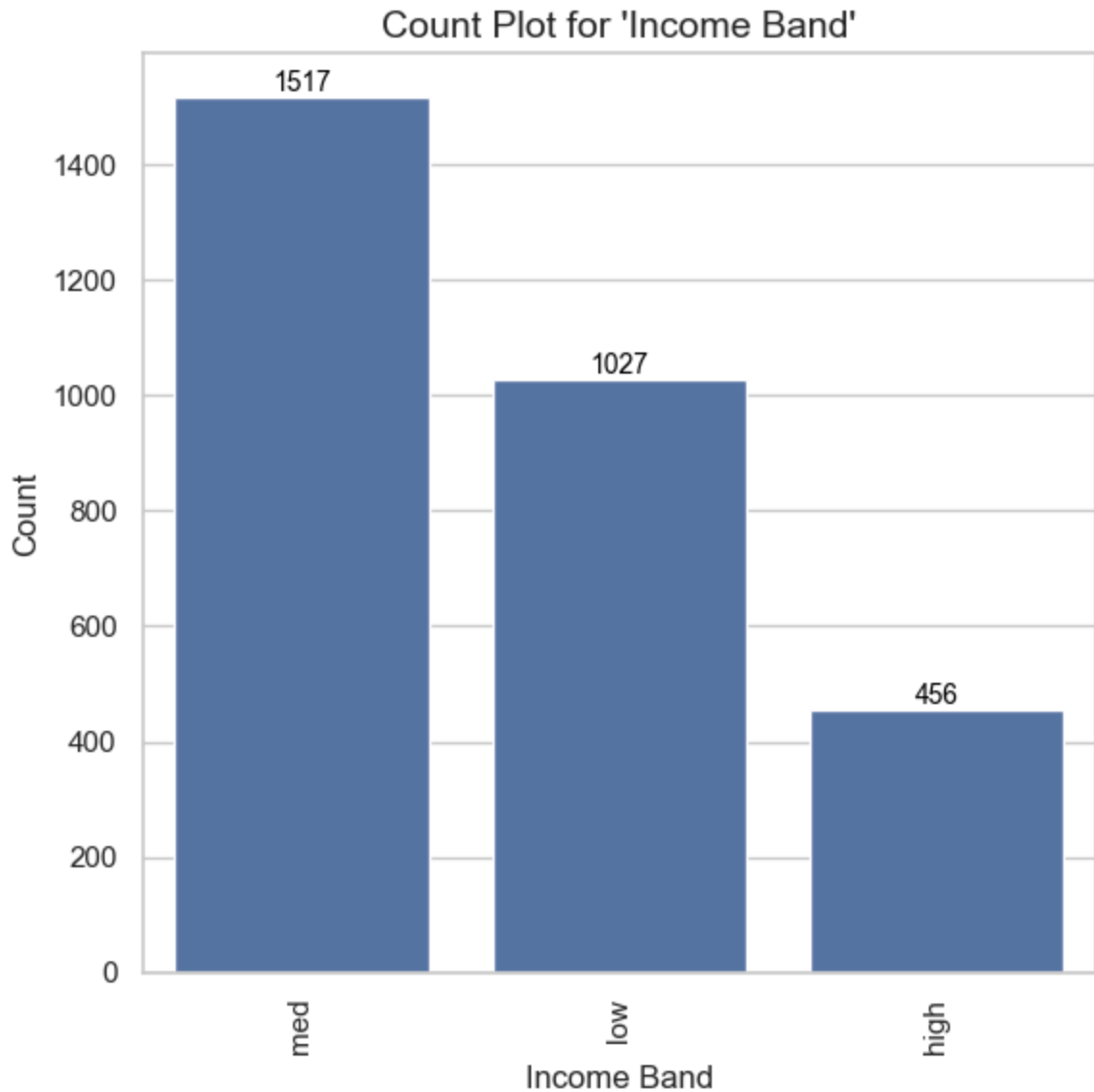
# Count Plot for 'Occupation'



Count Plot for 'Occupation' — bar chart showing Count (y-axis, 0 to ~28) versus Occupation (x-axis, with many overlapping labels).

Count Plot for 'Fee Structure'

Count Plot for 'Loyalty Classification'

Count Plot for 'Properties Owned'

Count Plot for 'Risk Weighting'

## Count Plot for 'Income Band'



# Bivariate Analysis

```
In [32]: sns.set(style="whitegrid")

for col in categorical_cols:
    if col == "GenderId":
        continue  # Skip GenderId since it's used as hue

    plt.figure(figsize=(5, 5))

    # Try to sort categories numerically if possible
    try:
        categories = sorted(df[col].dropna().unique(), key=lambda x: int(x))
    except:
        categories = df[col].value_counts().index

    # Plot with hue
    ax = sns.countplot(data=df, x=col, hue="Nationality", order=categories,
```

```python
# Add value labels on each bar
for p in ax.patches:
    height = int(p.get_height())
    if height > 0:
        ax.annotate(f'{height}',
                    (p.get_x() + p.get_width() / 2., height),
                    ha='center', va='bottom', fontsize=9, color='black',

# Formatting
plt.title(f"Bivariate Count Plot of '{col}' by GenderId", fontsize=14)
plt.xlabel(col, fontsize=12)
plt.ylabel("Count", fontsize=12)
plt.xticks(rotation=90)
plt.legend(title="GenderId")
plt.tight_layout()
plt.show()
```



Bivariate Count Plot of 'BRId' by GenderId

Bivariate Count Plot of 'IAId' by GenderId

# Bivariate Count Plot of 'Amount of Credit Cards' by GenderId
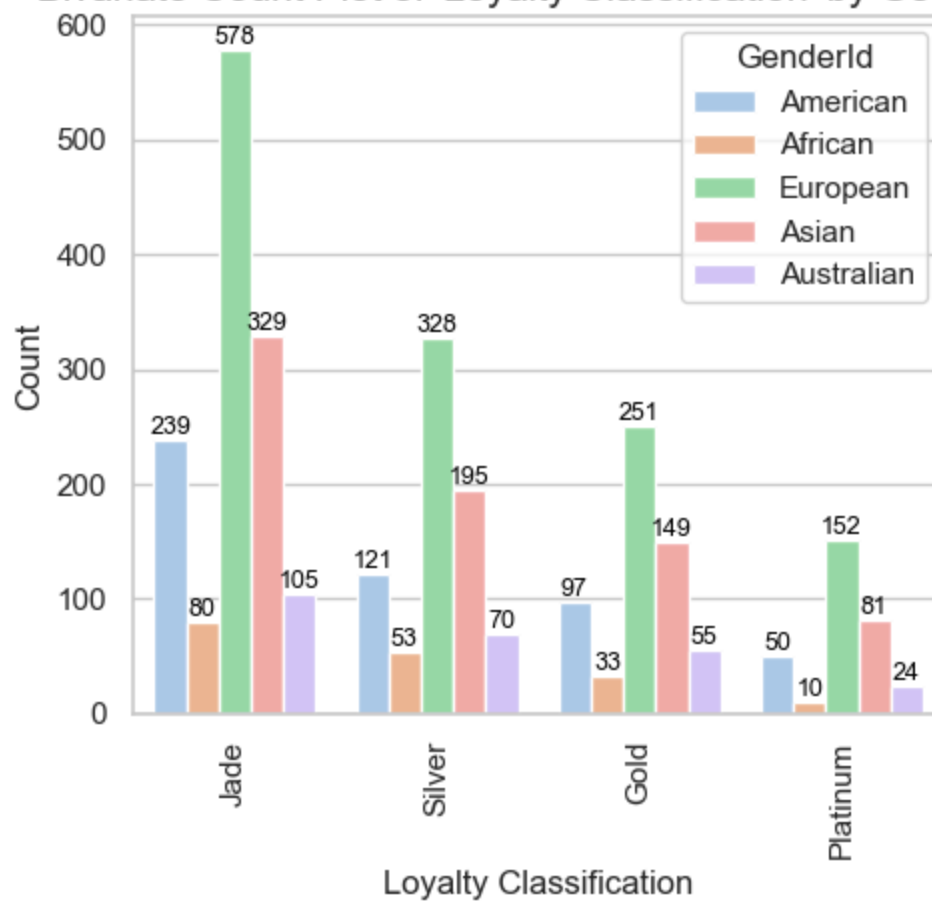
Bivariate Count Plot of 'Nationality' by GenderId

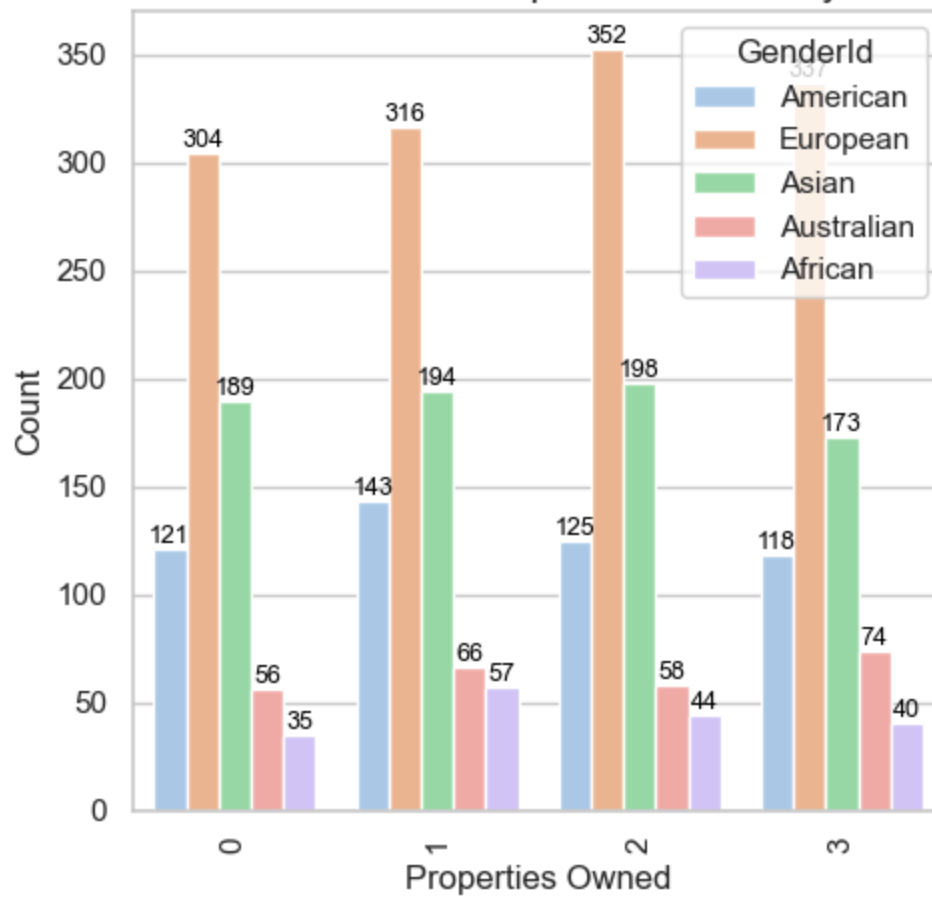Bivariate Count Plot of 'Occupation' by GenderId

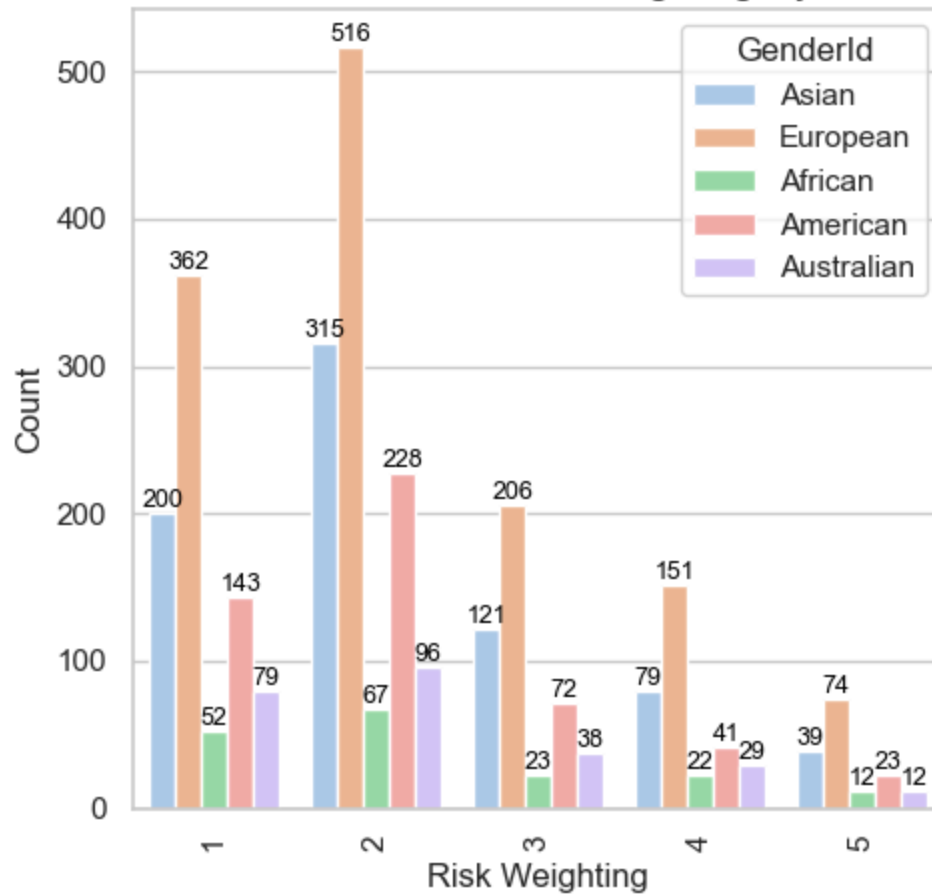Bivariate Count Plot of 'Fee Structure' by GenderId

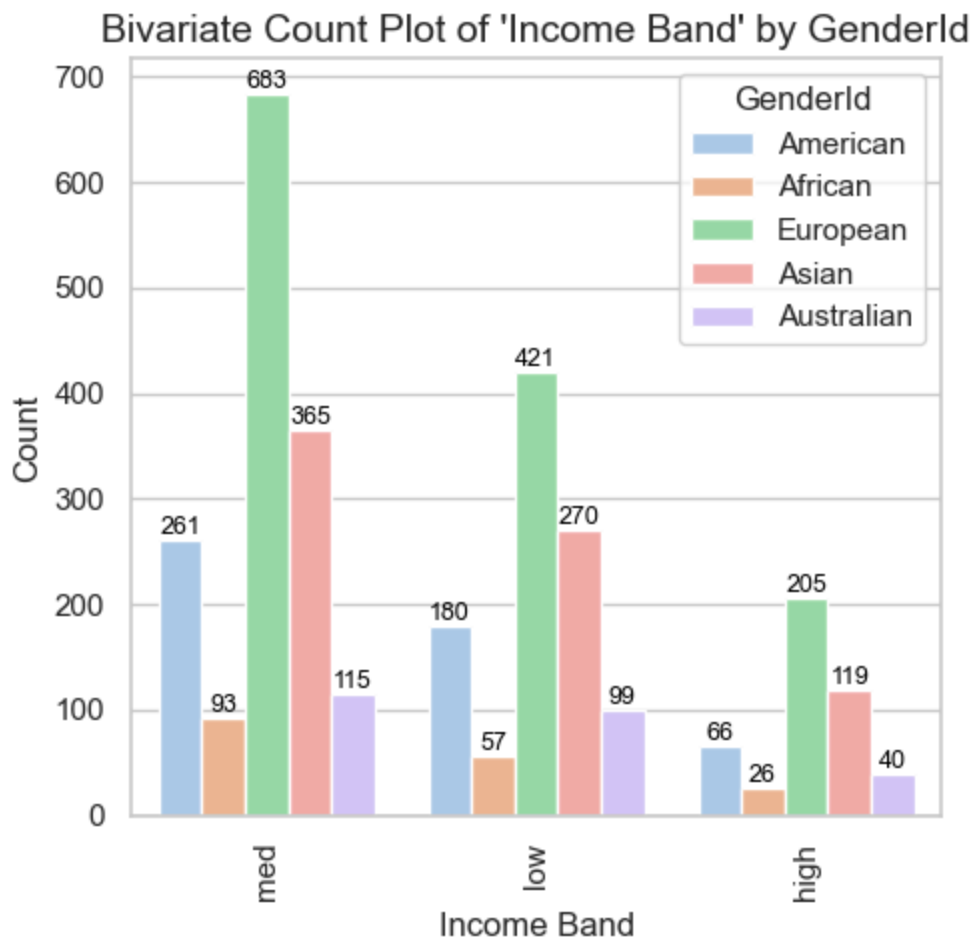Bivariate Count Plot of 'Loyalty Classification' by GenderId

Bivariate Count Plot of 'Properties Owned' by GenderId

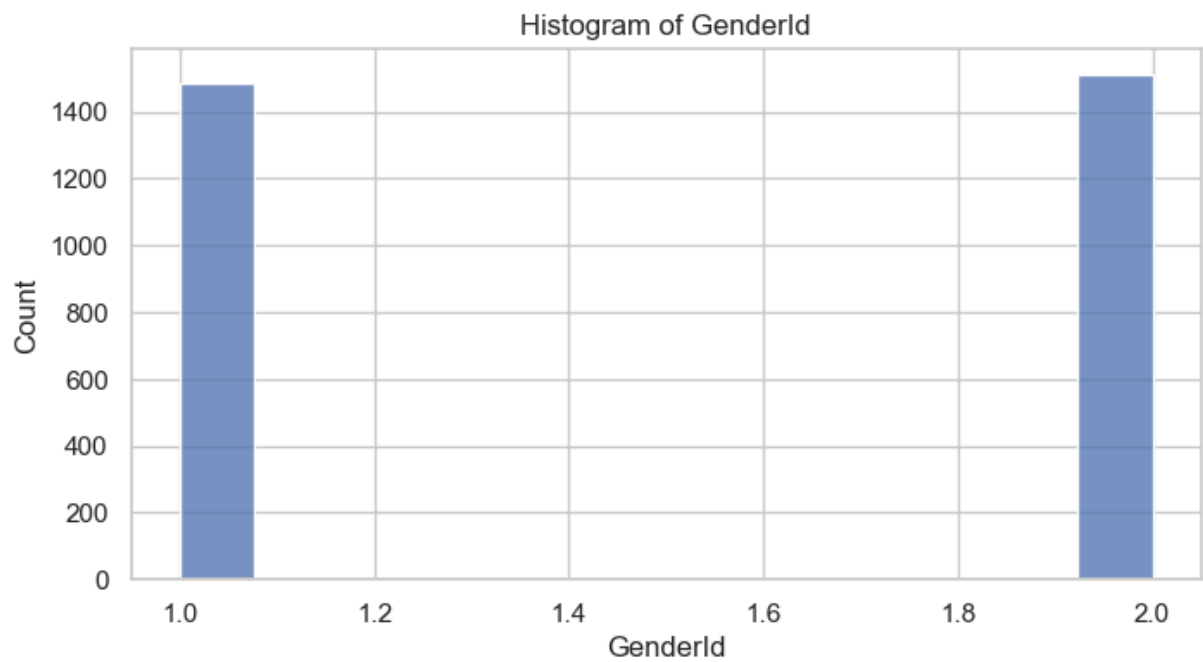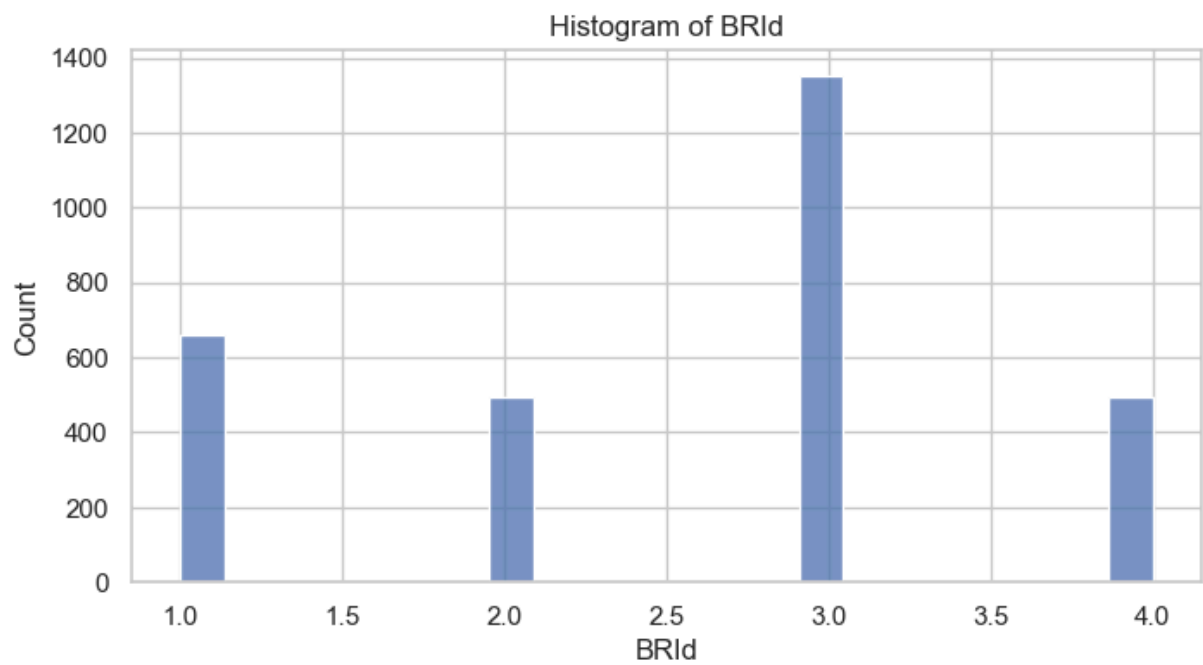Bivariate Count Plot of 'Risk Weighting' by GenderId

## Bivariate Count Plot of 'Income Band' by GenderId
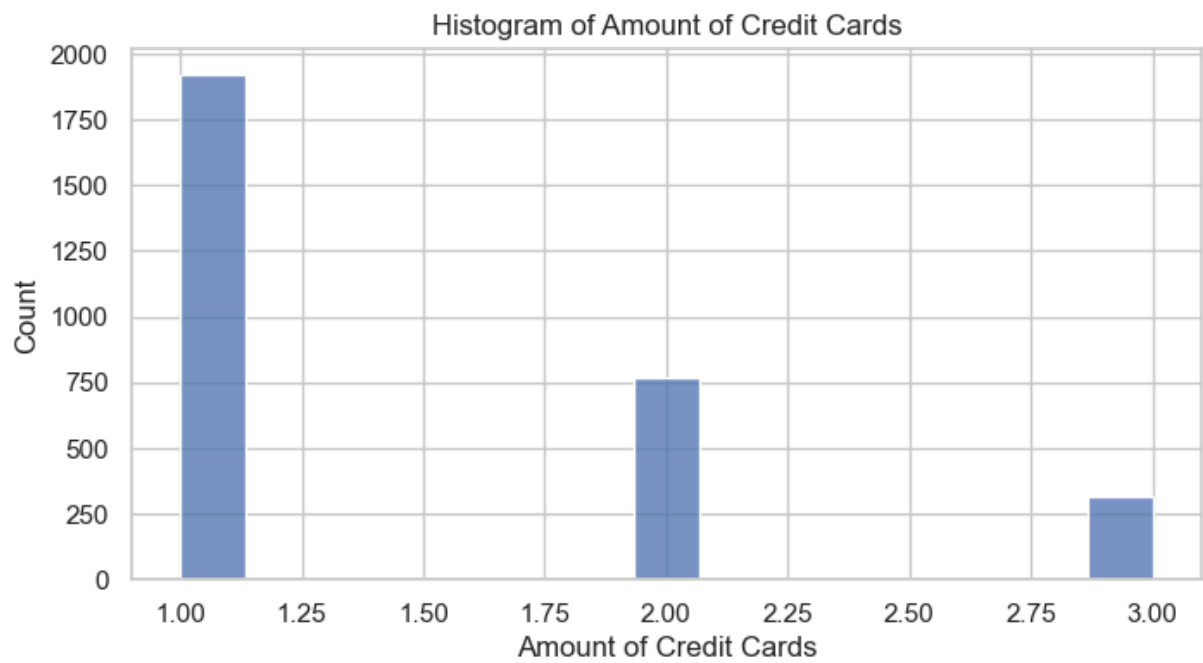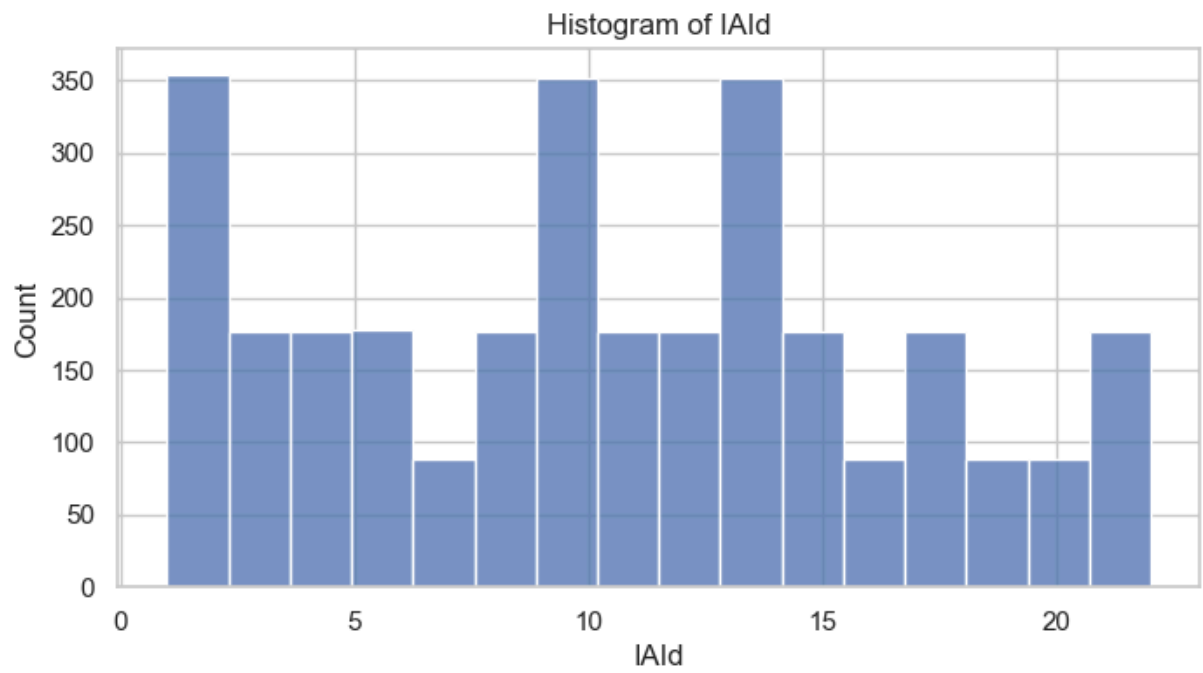


In [ ]: # Histplot of value counts for different occupation

In [38]:
```python
for col in categorical_cols:
    if col == "Occupation":
        continue  # Skip Occupation column

    plt.figure(figsize=(8, 4))
    sns.histplot(df[col])
    plt.title(f'Histogram of {col}')
    plt.xlabel(col)
    plt.ylabel("Count")
    plt.show()
```

Histogram of BRId



Histogram of GenderId

## Histogram of IAId



## Histogram of Amount of Credit Cards

Histogram of Nationality



Histogram of Fee Structure

Histogram of Loyalty Classification

Histogram of Properties Owned

## Histogram of Risk Weighting



## Histogram of Income Band



# Numerical Analysis

```
In [47]: numerical_cols=['Estimated Income', 'Superannuation Savings','Credit Card Ba

         #Univariate analysis and Visualization
         plt.figure(figsize=(10,4))
         for i,col in enumerate(numerical_cols):
             plt.subplot(4,3,i+1)
             sns.histplot(df[col],kde=True)
             plt.title(col)
             plt.show()
```

Estimated Income


Superannuation Savings


Credit Card Balance


Bank Loans


Bank Deposits

## Checking Accounts



## Saving Accounts



## Foreign Currency Account



## Business Lending



In [53]:
```python
numerical_cols=['Estimated Income', 'Superannuation Savings','Credit Card Ba

correlation_matrix= df[numerical_cols].corr()
plt.figure(figsize=(10,8))
sns.heatmap(correlation_matrix, annot=True, cmap='crest',fmt=".2f")
plt.title("Correlation matrix")
plt.show()
```

## Correlation matrix

| | Estimated Income | Superannuation Savings | Credit Card Balance | Bank Loans | Bank Deposits | Checking Accounts | Saving Accounts | Foreign Currency Account | Business Lending |
|---|---|---|---|---|---|---|---|---|---|
| Estimated Income | 1.00 | 0.37 | 0.30 | 0.33 | 0.26 | 0.29 | 0.26 | 0.31 | 0.33 |
| Superannuation Savings | 0.37 | 1.00 | 0.23 | 0.24 | 0.17 | 0.20 | 0.18 | 0.23 | 0.26 |
| Credit Card Balance | 0.30 | 0.23 | 1.00 | 0.37 | 0.38 | 0.30 | 0.28 | 0.36 | 0.35 |
| Bank Loans | 0.33 | 0.24 | 0.37 | 1.00 | 0.37 | 0.29 | 0.27 | 0.36 | 0.42 |
| Bank Deposits | 0.26 | 0.17 | 0.38 | 0.37 | 1.00 | 0.84 | 0.75 | 0.41 | 0.44 |
| Checking Accounts | 0.29 | 0.20 | 0.30 | 0.29 | 0.84 | 1.00 | 0.46 | 0.31 | 0.36 |
| Saving Accounts | 0.26 | 0.18 | 0.28 | 0.27 | 0.75 | 0.46 | 1.00 | 0.31 | 0.31 |
| Foreign Currency Account | 0.31 | 0.23 | 0.36 | 0.36 | 0.41 | 0.31 | 0.31 | 1.00 | 0.37 |
| Business Lending | 0.33 | 0.26 | 0.35 | 0.42 | 0.44 | 0.36 | 0.31 | 0.37 | 1.00 |

# Insight of EDA:

1. The strongest posituve correlation occurs among "Bank Deposits" and "Checking Accounts", "Saving Account" AND "Foreign Currency Account" imdicating customers wjo maintain high balances in one account type often hold substantial amount/funds across other accounts as well.