## Import libraries and read the file

```
In [1]:   1   # import libraries
          2
          3   import pandas as pd
          4   import seaborn as sns
          5   import numpy as np
          6   import matplotlib.pyplot as plt
          7   %matplotlib inline
```

```
In [2]:   1   # Read the census data
          2   census_data = pd.read_csv('census_data.csv')
```

## Study the dataset

```
In [3]:   1   #Check the number of rows and columns
          2   census_data.shape
```

```
Out[3]:  (8646, 11)
```

```
In [4]:   1  # print out the first five rows to have an overview of the dataset
          2  census_data.head()
```

Out[4]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | George Avenue | Harry | James | 60 | Head | Single | Male | Unemployed | None | None |
| 1 | 2 | George Avenue | Anne | Johnson | 34 | Head | Married | Female | Corporate treasurer | None | None |
| 2 | 2 | George Avenue | Jack | Johnson | 36 | Husband | Married | Male | Product/process development scientist | None | None |
| 3 | 2 | George Avenue | Guy | Johnson | 12 | Son | NaN | Male | Student | None | NaN |
| 4 | 3 | George Avenue | Simon | Smith | 79 | Head | Single | Male | Retired Tour manager | Physical Disability | Jewish |

```
In [5]:   1  # Displays the data type, non null count and number of data entry
          2  census_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8646 entries, 0 to 8645
Data columns (total 11 columns):
 #   Column                         Non-Null Count  Dtype
---  ------                         --------------  -----
 0   House Number                   8646 non-null   int64
 1   Street                         8646 non-null   object
 2   First Name                     8646 non-null   object
 3   Surname                        8646 non-null   object
 4   Age                            8646 non-null   object
 5   Relationship to Head of House  8646 non-null   object
 6   Marital Status                 6419 non-null   object
 7   Gender                         8646 non-null   object
 8   Occupation                     8646 non-null   object
 9   Infirmity                      8646 non-null   object
 10  Religion                       6373 non-null   object
dtypes: int64(1), object(10)
memory usage: 743.1+ KB
```

```python
In [6]:   1  # Check the total number of missing values
          2  census_data.isna().sum()
```

```
Out[6]:  House Number                     0
         Street                           0
         First Name                       0
         Surname                          0
         Age                              0
         Relationship to Head of House    0
         Marital Status                2227
         Gender                           0
         Occupation                       0
         Infirmity                        0
         Religion                      2273
         dtype: int64
```

```python
In [7]:   1  # Check for duplicate
          2  duplicate = census_data.duplicated()
          3  census_data[duplicate]
```

Out[7]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **7309** | 1 | Leedsbox Crescent | Ashleigh | Osborne | 15 | Daughter | NaN | Female | Student | None | NaN |

```python
In [8]:   1  # Drop the duplicate
          2  census_data = census_data.drop(7309)
```

# Data Cleaning

## Age

```
In [9]:   1  # Check the unique entries for age
          2  print(census_data['Age'].unique())
```

```
['60' '34' '36' '12' '79' '35' '61' '24' '3' '75' '52' '14' '11' '42' '25'
 '28' '40' '57' '55' '22' '18' '43' '51' '0' '21' '45' '17' '16' '13' '9'
 '65' '32' '31' '8' '56' '39' '7' '41' '27' '78' '30' '29' '15' '54' '19'
 '84' '38' '33' '6' '1' '48' '10' '5' '49' '46' '26' '50' '53' '63' '4'
 '44' '47' '2' '23' '64' '37' '58' '66' '67' '71' '72' '20' '62' '68' '73'
 '74' '69' '81' '70' '59' '89' '105' '87' '80' '77' '76' ' ' '82' '88'
 '49.16040882016717' '54.16040882016717' '3.0' '85' '99' '101' '83'
 '69.13036593215614' '67.13036593215614' '103' '90' '93' '86' '96'
 '85.66111048772531' '87.66111048772531' '34.0' '30.0' '26.0' '91' '102'
 '83.52432893335205' '26.999999999999993' '23.999999999999993'
 '21.999999999999993' '16.999999999999993' '92' '97' '69.13473801820774'
 '15.000000000000007' '13.000000000000007' '10.000000000000007' '98'
 '50.53760781824045' '53.53760781824045' '0.0']
```

```
In [10]:  1  # Check for blank entries in Age
          2  census_data[census_data['Age'] == ' ']
```

Out[10]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **460** | 18 | Smith Gateway | Dominic | Griffiths | | Son | NaN | Male | Student | None | NaN |

```
In [11]:  1  # Drops the line, since it is just a row, and we cannot predict the age of this person
          2  census_data = census_data.drop(460)
```

In [12]:
```python
# Convert the Age column to integer
census_data['Age'] = census_data['Age'].astype(float).round(0).astype(int)
```
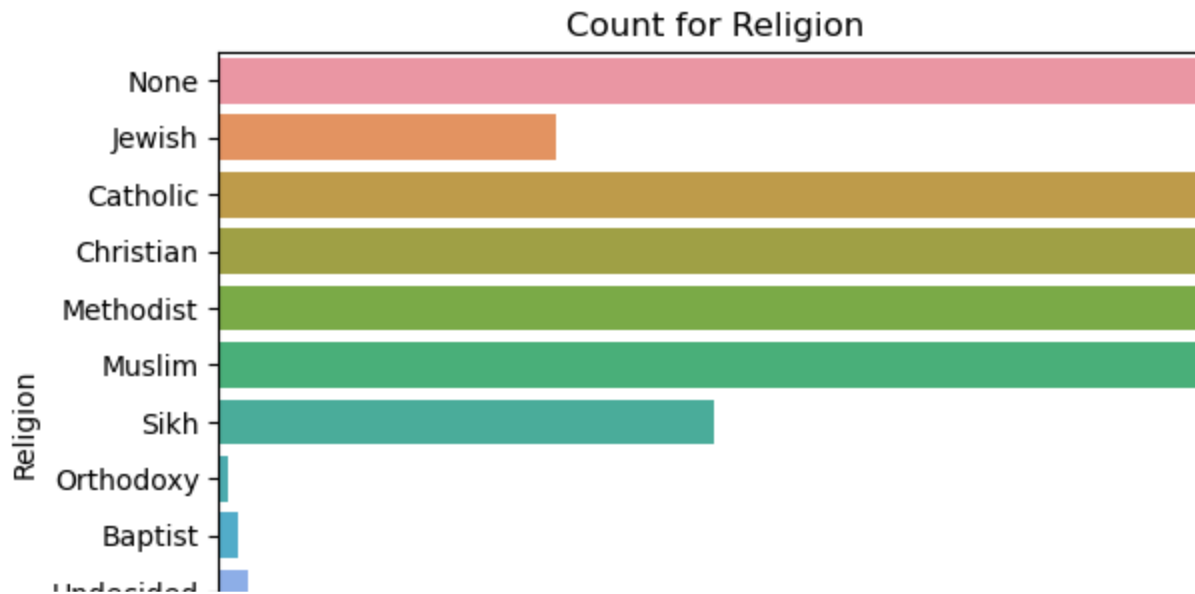
## Religion

In [13]:
```python
# Check the unique entries in Religion
print(census_data['Religion'].unique())
```

```
['None' nan 'Jewish' 'Catholic' 'Christian' 'Methodist' 'Muslim' 'Sikh'
 'Orthodoxy' 'Baptist' 'Undecided' 'Buddist' ' ' 'Nope']
```

```
In [14]:  1  # Countplot for religion to visualise the data
          2  religion_bar = sns.countplot(census_data, y = census_data['Religion'])
          3
          4  # sets the x-axis limit from 0-100, to easily visualised the entries with small values.
          5  plt.xlim(0,100)
          6
          7  # Title of the plot
          8  religion_bar.set(title = 'Count for Religion')
```

Out[14]: [Text(0.5, 1.0, 'Count for Religion')]

```
In [15]:  1  # Check for blank entries in Religion
          2  census_data[census_data['Religion'] == ' ']
```

Out[15]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **7069** | 4 | Parsons Stream | Neil | Hall | 34 | Husband | Married | Male | Mining engineer | None | |
| **7809** | 1 | Cox Drive | Valerie | Arnold | 57 | Head | Single | Female | Unemployed | None | |
| **8385** | 57 | George Lane | Debra | Davies | 31 | Head | Married | Female | Barista | None | |
| **8433** | 67 | George Lane | Ashleigh | Martin | 38 | Lodger | Single | Female | Minerals surveyor | None | |

```
In [16]:  1  # Check the range of each entry to find a family tie
          2  census_data[7806:7810]
```

Out[16]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **7808** | 4 | Kelly Mountain | Denise | Thompson | 10 | Daughter | NaN | Female | Student | None | NaN |
| **7809** | 1 | Cox Drive | Valerie | Arnold | 57 | Head | Single | Female | Unemployed | None | |
| **7810** | 1 | Cox Drive | Katie | Arnold | 23 | Daughter | Single | Female | Health promotion specialist | None | Christian |
| **7811** | 2 | Cox Drive | Kyle | Perkins | 72 | Head | Widowed | Male | Retired Film/video editor | None | Christian |

```
In [17]:  1  # Replace row 7809 with Christian since the daughter in row 7810 is a Christian
          2  census_data.loc[7809,'Religion'] = 'Christian'
```

```
In [18]:  1  # Replace the remaining blanks to 'None' since they are adult
          2  census_data['Religion'].replace(' ', 'None', regex = False, inplace = True)
```

```
In [19]:   1  # Replace the entry 'Nope' to 'None' for consistency
           2  census_data['Religion'].replace('Nope', 'None', regex = True, inplace = True)
           3
           4  # Replace the entry 'Buddist' to 'Buddhist' for consistency
           5  census_data['Religion'].replace('Buddist', 'Buddhist', regex = True, inplace = True)
```

**Replace Methodist, Baptist and Orthodoxy to Christian, as these are not religion**

```
In [20]:   1  # Replace Methodist, Baptist and Orthodoxy to Cristian
           2  census_data['Religion'].replace('Methodist', 'Christian', regex = True, inplace = True)
           3  census_data['Religion'].replace('Baptist', 'Christian', regex = True, inplace = True)
           4  census_data['Religion'].replace('Orthodoxy', 'Christian', regex = True, inplace = True)
```

**Fill the missing value of religion with the 'Head' of house, since children can adopt their parents religion**

```
In [21]:   1  # Filter age less than 18 from religion missing values
           2  religion_age_data = census_data[(census_data['Religion'].isna()) & (census_data['Age'] < 18)]
           3
           4  '# Iterate through the data
           5  for index, row in religion_age_data.iterrows():
           6      # Find the row of each head of house
           7      head_of_house = census_data[(census_data['House Number'] == row['House Number']) &
           8                                  (census_data['Street'] == row['Street']) &
           9                                  (census_data['Relationship to Head of House'] == 'Head')]
          10
          11      # If head of house is not empty, replace the missing values with it
          12      if not head_of_house.empty:
          13          census_data.loc[index, 'Religion'] = head_of_house.iloc[0]['Religion']
```

```
  File "C:\Users\rosem\AppData\Local\Temp\ipykernel_11532\571342830.py", line 4
    '# Iterate through the data
                               ^
SyntaxError: EOL while scanning string literal
```

```
In [22]:  1  # Replace other missing values with None
          2  census_data['Religion'] = census_data['Religion'].fillna('None')
```
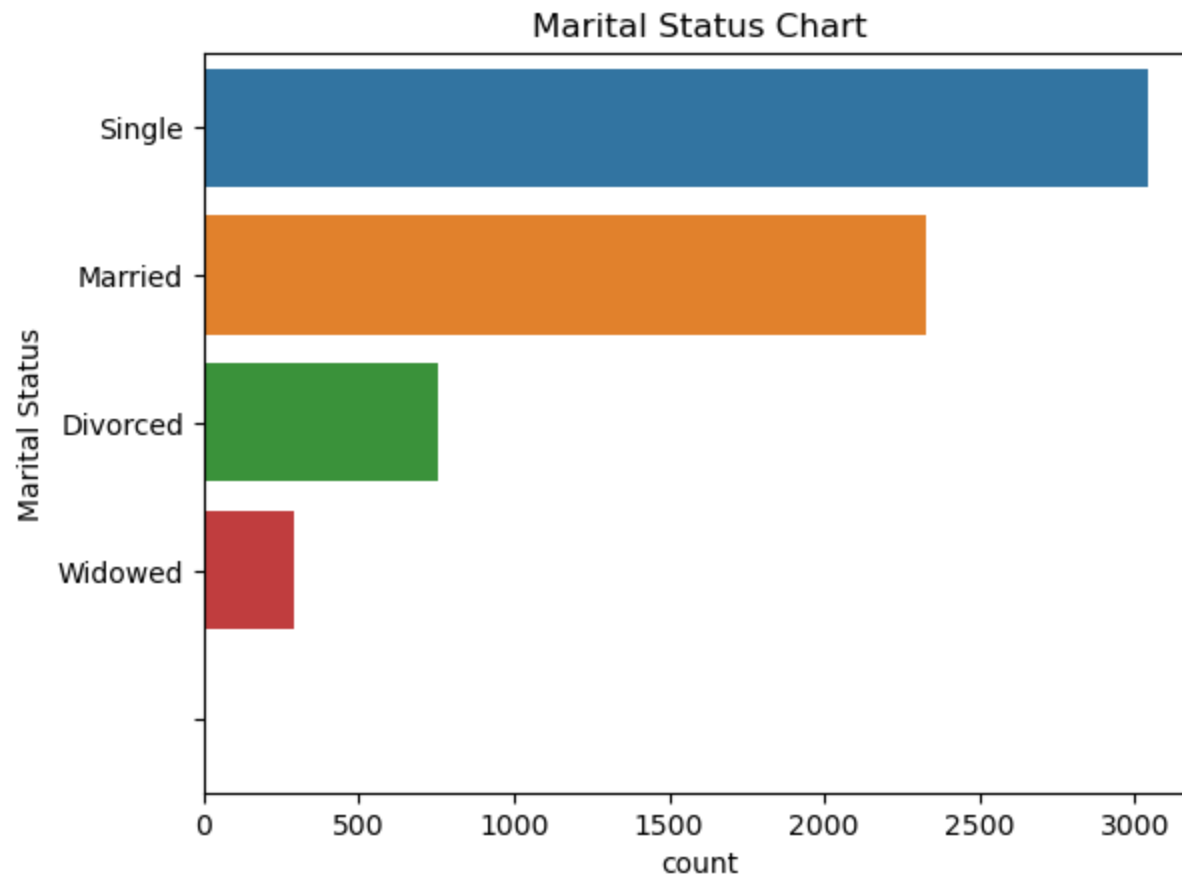
## Marital Status

```
In [23]:  1  # print the unique entries in Marital Status
          2  print(census_data['Marital Status'].unique())
```

['Single' 'Married' nan 'Divorced' 'Widowed' ' ']

```
In [24]:   1  # Plot a chart for Marital Status
           2  marital_status_count = sns.countplot(census_data, y = census_data['Marital Status'])
           3
           4  # Title the plot
           5  marital_status_count.set(title='Marital Status Chart')
```

Out[24]: [Text(0.5, 1.0, 'Marital Status Chart')]

```
In [25]:    1  # Check blank entries for Marital Status
            2  census_data[census_data['Marital Status'] == ' ']
```

Out[25]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **3205** | 43 | Morgan Fords | Diana | Robinson | 39 | Wife | | Female | Hospital pharmacist | None | None |

```
In [26]:    1  # Check the range of each entry to find a family tie
            2  census_data[3203:3209]
```

Out[26]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **3204** | 43 | Morgan Fords | Peter | Robinson | 42 | Head | Married | Male | Designer, interior/spatial | None | Catholic |
| **3205** | 43 | Morgan Fords | Diana | Robinson | 39 | Wife | | Female | Hospital pharmacist | None | None |
| **3206** | 43 | Morgan Fords | Wayne | Robinson | 12 | Son | NaN | Male | Student | None | None |
| **3207** | 43 | Morgan Fords | Dale | Robinson | 5 | Son | NaN | Male | Student | None | None |
| **3208** | 43 | Morgan Fords | Charles | Robinson | 3 | Son | NaN | Male | Child | None | None |
| **3209** | 44 | Morgan Fords | Beverley | Williams | 39 | Head | Divorced | Female | Copywriter, advertising | None | None |

```
In [27]:    1  # Replace row 3205 with Married, since the husband who is the head of house is Married
            2  census_data.loc[3205,'Marital Status'] = 'Married'
```

**Change children marital status to 'Not Available' because they do not have a marital status yet until they are 18 years old**

```
In [28]:   1  # Find missing values in marital status for children
           2  children_missing_value = (census_data['Age'] < 18) & (census_data['Marital Status'].isna())
           3
           4  # Replace missing values with 'Not Available'
           5  census_data.loc[children_missing_value, 'Marital Status'] = 'Not Available'
```

## Gender

```
In [29]:   1  census_data['Gender'].unique()
```

Out[29]: array(['Male', 'Female', ' '], dtype=object)

```
In [30]:   1  # checks for blank entries in Gender
           2  census_data[census_data['Gender'] == ' ']
```

Out[30]:

|  | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **6013** | 9 | Lime Street | Elizabeth | Dobson | 4 | Daughter | Not Available | | Child | None | None |
| **7538** | 37 | Leedsbox Crescent | Liam | Yates | 66 | Head | Married | | Television production assistant | None | Christian |

```
In [31]:   1  # Check the range of each entry to find a family tie
           2  census_data[7535:7540]
```

Out[31]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **7537** | 36 | Leedsbox Crescent | Danielle | Palmer | 69 | Head | Widowed | Female | Retired Engineer, mining | None | Christian |
| **7538** | 37 | Leedsbox Crescent | Liam | Yates | 66 | Head | Married | | Television production assistant | None | Christian |
| **7539** | 37 | Leedsbox Crescent | Hayley | Yates | 66 | Wife | Married | Female | Solicitor, Scotland | None | Christian |
| **7540** | 37 | Leedsbox Crescent | Dylan | Yates | 41 | Son | Single | Male | Materials engineer | None | None |
| **7541** | 37 | Leedsbox Crescent | Terry | Yates | 38 | Son | Single | Male | Designer, fashion/clothing | None | Christian |

```
In [32]:   1  # Replace with Female, since her status is 'Daughter'
           2  census_data.loc[6013,'Gender'] = 'Female'
```

```
In [33]:   1  # Replace with Male, since he is the head and has a wife
           2  census_data.loc[7538,'Gender'] = 'Male'
```
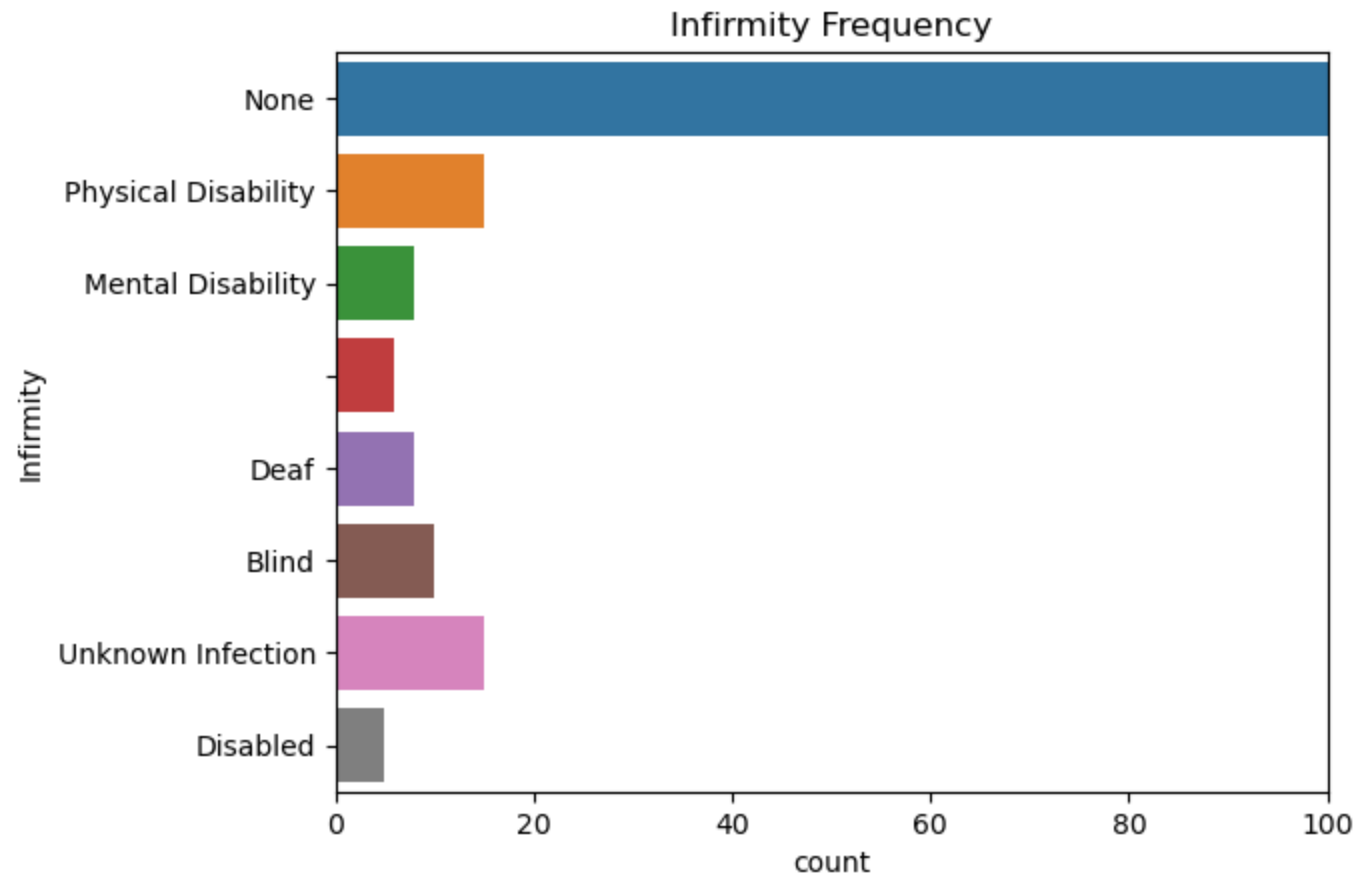
## Infirmity

```
In [34]:   1  # Check for unique entries in Infirmity
           2  census_data['Infirmity'].unique()
```

Out[34]: array(['None', 'Physical Disability', 'Mental Disability', ' ', 'Deaf',
               'Blind', 'Unknown Infection', 'Disabled'], dtype=object)

```
In [35]:  1  # Count plot for infirmity
          2  infirmity_count = sns.countplot(census_data, y = census_data['Infirmity'])
          3
          4  # Limit the count to 100
          5  plt.xlim(0,100)
          6
          7  # Title the plot
          8  infirmity_count.set(title='Infirmity Frequency')
```

Out[35]: [Text(0.5, 1.0, 'Infirmity Frequency')]

```
In [36]:   1  # Check for blank entries in Infirmity
           2  census_data[census_data['Infirmity']== ' ']
```

Out[36]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **556** | 1 | Morgan View | Sean | Howe | 15 | Son | Not Available | Male | Student | | None |
| **909** | 15 | Newfound Station | Lynda | Murphy | 24 | Head | Single | Female | Public affairs consultant | | Christian |
| **1120** | 24 | Palmer Crescent | Garry | Burns | 52 | Husband | Married | Male | Actuary | | None |
| **4244** | 47 | Madridgate Drive | Fiona | Lloyd | 79 | Wife | Married | Female | Retired Contractor | | Christian |
| **6047** | 12 | Graham Road | Caroline | Bruce | 46 | Head | Divorced | Female | Chartered management accountant | | Catholic |
| **7727** | 9 | Salmon Lane | Holly | Francis | 40 | Head | Single | Female | Programmer, systems | | Catholic |

```
In [37]:   1  # Replace all blank entries with None
           2  census_data['Infirmity']= census_data['Infirmity'].replace(' ', 'None')
```

## Surname

```
In [38]:   1  # Check for blank entries in Surname
           2  census_data[census_data['Surname']== ' ']
```

Out[38]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **2123** | 24 | Morley Lodge | Simon | | 56 | None | Single | Male | Information officer | None | Christian |
| **3168** | 33 | Morgan Fords | Stephanie | | 8 | Daughter | Not Available | Female | Student | None | None |
```

```
In [39]:   1  # Check the range to identify family relationship
           2  census_data[2121:2126]
           3
           4  # There is no family member associated with 2123, I will ignore the data and continue.
           5  # A blank surname will not affect our analysis
```

Out[39]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **2122** | 24 | Morley Lodge | Caroline | Barber | 38 | None | Single | Female | Pharmacologist | None | Christian |
| **2123** | 24 | Morley Lodge | Simon | | 56 | None | Single | Male | Information officer | None | Christian |
| **2124** | 25 | Morley Lodge | Dylan | Griffiths | 34 | Head | Single | Male | Financial manager | None | Catholic |
| **2125** | 25 | Morley Lodge | Eleanor | Griffiths | 43 | Cousin | Single | Female | Teacher, secondary school | None | Catholic |
| **2126** | 26 | Morley Lodge | Lindsey | Smith | 44 | Head | Single | Female | Unemployed | None | Catholic |

```
In [40]:   1  # Check the range to identify family relationship
           2  census_data[3165:3170]
```

Out[40]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **3166** | 33 | Morgan Fords | Lorraine | Griffin | 31 | Head | Married | Female | Unemployed | None | None |
| **3167** | 33 | Morgan Fords | Henry | Griffin | 31 | Husband | Married | Male | Acupuncturist | None | None |
| **3168** | 33 | Morgan Fords | Stephanie | | 8 | Daughter | Not Available | Female | Student | None | None |
| **3169** | 33 | Morgan Fords | Francis | Griffin | 4 | Son | Not Available | Male | Child | None | None |
| **3170** | 33 | Morgan Fords | Kathryn | Dobson | 2 | Daughter | Not Available | Female | Child | None | None |

```
In [41]:   1  # Input surname for 3168 with Griffin since they are family
           2  census_data.loc[3168, 'Surname']='Griffin'
```

## First Name

```
In [42]:   1  # Check for blank cells in First Name
           2  census_data[census_data['First Name']== ' ']
           3
           4  #First Name is a unique value, I will ignore it and continue, as this will not affect our further anlysis
```

Out[42]:

| | House Number | Street | First Name | Surname | Age | Relationship to Head of House | Marital Status | Gender | Occupation | Infirmity | Religion |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **618** | 19 | Morgan View | | Ali | 8 | Son | Not Available | Male | Student | None | None |
| **3266** | 4 | Simmons Course | | Wong | 9 | Son | Not Available | Male | Student | None | None |
| **3916** | 6 | ExcaliburBells Road | | Doyle | 5 | Son | Not Available | Male | Student | None | None |

## House Number, Street

```
In [43]:   1  len(census_data[census_data['House Number']== ' '])
```

Out[43]:  0

```
In [44]:   1  len(census_data[census_data['Street'] ==' '])
```

Out[44]:  0

```
In [45]:   1  len(census_data[census_data['Relationship to Head of House']== ' '])
```

Out[45]:  0

```
In [46]:   1  len(census_data[census_data['Occupation']== ' '])
```

Out[46]:  0

## Discussion / Analysis

```
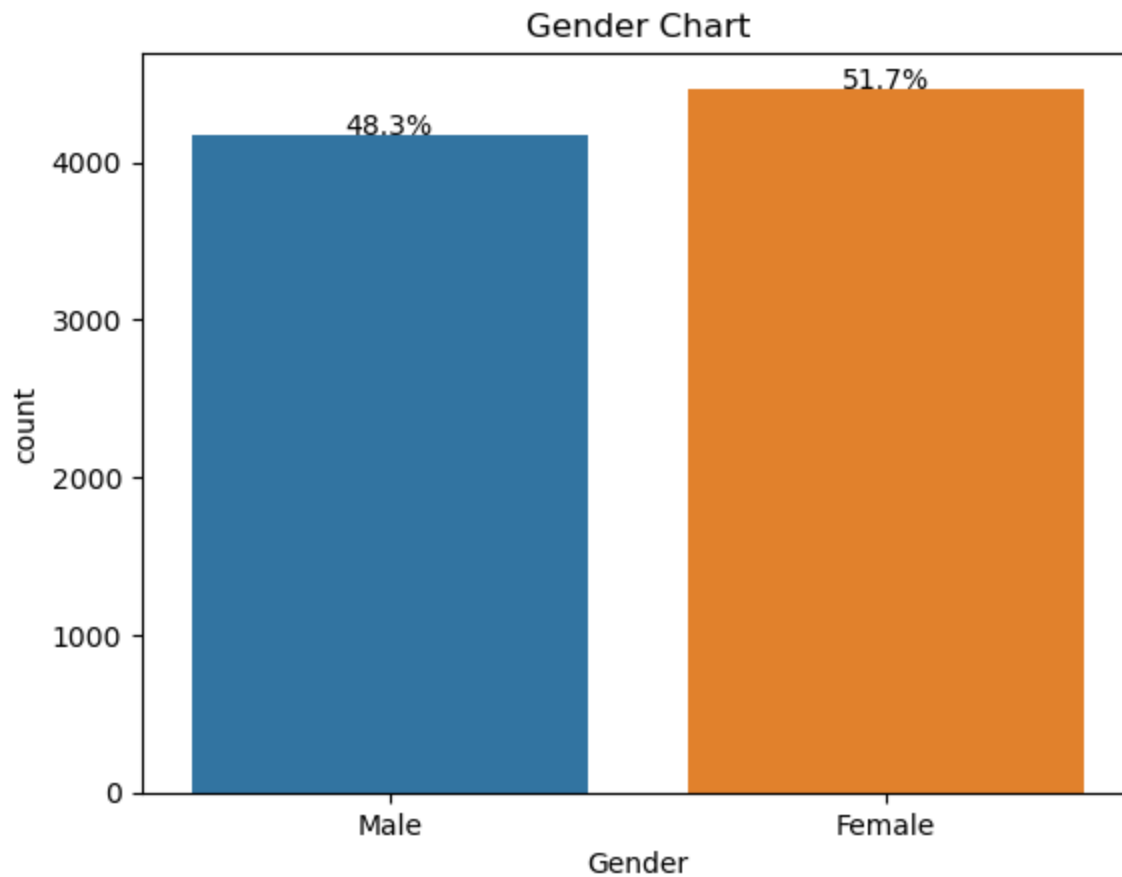In [47]:   1  # Save the cleaned data to a csv file
           2  census_data.to_csv('cleaned_data.csv', index = False)
```

```
In [48]:   1  # Read the cleaned csv file
           2  new_data = pd.read_csv('cleaned_data.csv')
```

## Age Distribution

```python
# Create a plot for the Gender column
gender_plot = sns.countplot(data=new_data, x='Gender')

# Calculate the total count of gender
total_frequency = len(new_data['Gender'])

# Loop through each bar in the plot
for p in gender_plot.patches:
    # Get the percentage value for the bar
    percentage = p.get_height() / total_frequency * 100

    # Add the percentage as text above the bar
    gender_plot.annotate(f'{percentage:.1f}%', (p.get_x() + p.get_width() / 2., p.get_height()), ha='center', va=

# Give the plot a title and name the x-axis
gender_plot.set_title('Gender Chart')
gender_plot.set_xlabel('Gender')
```

Out[49]: Text(0.5, 0, 'Gender')

## Gender Chart



## Population Pyramid

```
In [50]:   1  # Create a list of age group with bin width 10, and add it to the existing new_data
           2
           3  bin_width = [0,9,19,29,39,49,59,69,79,89,99,105]
           4  new_data['Age Group'] = pd.cut(new_data['Age'],bin_width,
           5                           labels = ['0-9','10-19','20-29','30-39','40-49','50-59','60-69','70-79','80-89','9
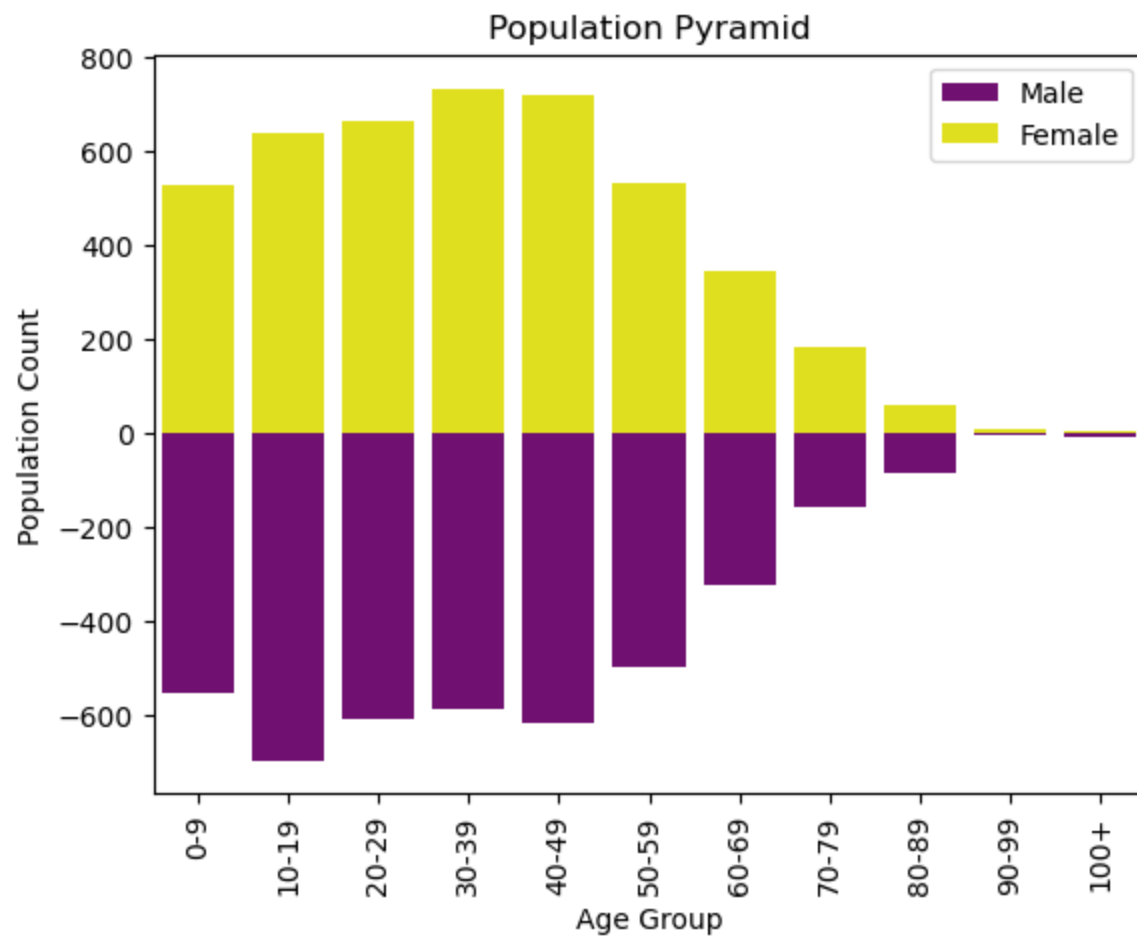```

```python
In [51]:   1  # Group age by gender
           2  group_age_bygender = new_data.groupby(["Age Group", "Gender"]).size().unstack()
```

```python
In [52]:   1  # Create a variable name for male age after multiplying it by 1
           2  male_age = group_age_bygender["Male"] * -1
           3
           4  # Create a variable name for female age
           5  female_age = group_age_bygender["Female"]
           6
           7  # Use female data as the age group
           8  age_group = group_age_bygender.index
```

```python
In [53]:   1  # Create the dataframe for Age
           2  age_df = pd.DataFrame({'Age': age_group,
           3  'Male': male_age,
           4  'Female': female_age})
```

```
In [54]:   1  # Create a barplot for the male population
           2  age_pyramid = sns.barplot(y='Male', x='Age', data=age_df,  color=('purple'), label='Male')
           3
           4  # Create a second barplot on top of the first one for the female population
           5  age_pyramid = sns.barplot(y='Female', x='Age', data=age_df,  color=('yellow'), label='Female')
           6
           7  # Add a legend to the plot to differentiate between male and female bars
           8  age_pyramid.legend()
           9
          10  # Rotate x-labels by 90 degrees
          11  age_pyramid.set_xticklabels(age_pyramid.get_xticklabels(), rotation=90)
          12
          13  # Title the plot, x and y axis
          14  age_pyramid.set(ylabel='Population Count', xlabel='Age Group')
          15  plt.title('Population Pyramid')
```

Out[54]:  Text(0.5, 1.0, 'Population Pyramid')

Population Pyramid

```
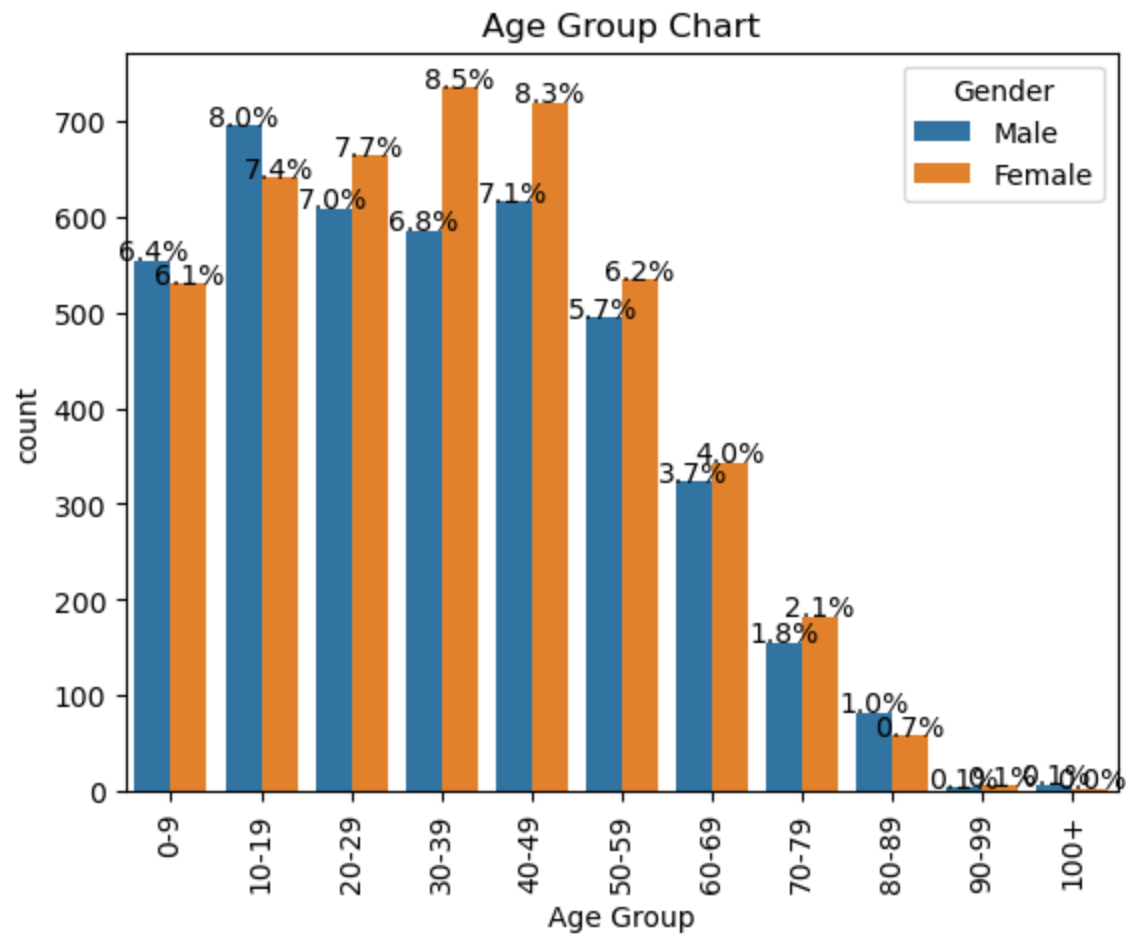In [55]:    1  # Violin plot for Age distribution
            2  violin_age_plot = sns.violinplot(new_data, x = 'Age')
            3
            4  # Title the plot
            5  violin_age_plot.set(title= 'Age Distribution Chart')
```

Out[55]:  [Text(0.5, 1.0, 'Age Distribution Chart')]



Age Distribution Chart

```
1  # Create a plot for the Gender column
2  age_group_plot = sns.countplot(data=new_data, x='Age Group', hue = 'Gender')
3
4  # Calculate the total count of age group
5  total_frequency = len(new_data['Age Group'])
6
7  # Loop through each bar in the plot
8  for p in age_group_plot.patches:
9      # Get the percentage value for the bar
10     percentage = p.get_height() / total_frequency * 100
11
12     # Add the percentage as text above the bar
13     age_group_plot.annotate(f'{percentage:.1f}%', (p.get_x() + p.get_width() / 2., p.get_height()), ha='center',
14
15 # Rotate the x-axis label
16 age_group_plot.set_xticklabels(age_group_plot.get_xticklabels(), rotation=90)
17
18 # Give the plot a title and name the x-axis
19 age_group_plot.set(title = 'Age Group Chart')
20 age_group_plot.set_xlabel('Age Group')
```

Text(0.5, 0, 'Age Group')

Age Group Chart

**Occupation Level**

```
In [57]:    1  def occupation_status(occupation):
            2
            3      if 'University Student'in occupation:
            4          return 'University Student'
            5      if 'PhD Student'in occupation:
            6          return 'PhD Student'
            7      if 'Student' in occupation:
            8          return 'Students'
            9      elif 'Retired' in occupation:
           10          return 'Retired'
           11      elif 'Child' in occupation and not ("Child psychotherapist" in occupation or "Nurse,children's" in occupation
           12          return 'Child'
           13      elif 'Unemployed' in occupation:
           14          return 'Unemployed'
           15      else:
           16          return 'Employed'
```

```
In [58]:    1  # Create a new column called Occupation Status
            2  new_data['Occupation Status']= new_data['Occupation'].apply(occupation_status)
```

```
In [59]:  1  # Create a plot for the Gender column
          2  occu_plot = sns.countplot(data=new_data, x='Occupation Status', hue = 'Gender')
          3
          4  # Calculate the total count of age group
          5  total_frequency = len(new_data['Age Group'])
          6
          7  # Loop through each bar in the plot
          8  for p in occu_plot.patches:
          9      # Get the percentage value for the bar
         10      percentage = p.get_height() / total_frequency * 100
         11
         12      # Add the percentage as text above the bar
         13      occu_plot.annotate(f'{percentage:.1f}%', (p.get_x() + p.get_width() / 2., p.get_height()), ha='center', va='b
         14
         15
         16  # Rotate the x-axis label
         17  occu_plot.set_xticklabels(occu_plot.get_xticklabels(), rotation=45)
         18
         19  # Give the plot a title and name the x-axis
         20  occu_plot.set(title = 'Occupation Status Plot')
         21  occu_plot.set_xlabel('Occupation Status')
```

Out[59]: Text(0.5, 0, 'Occupation Status')

```
In [60]:   1  # Create variables for employed
           2  employed = new_data[new_data['Occupation Status'] == 'Employed']['Occupation Status'].count()
           3
           4  # Create variables for unemployed
           5  unemployed = new_data[new_data['Occupation Status'] == 'Unemployed']['Occupation Status'].count()
           6
           7  # Create variables for workforce_age
           8  workforce_age = new_data[(new_data['Age'] >= 18) & (new_data['Age'] <= 65)]['Age'].count()
```

## Workforce Calculation

```
In [61]:   1  # Calculate the total population
           2  total_population = new_data.shape[0]
           3
           4  print(f'The total population is {total_population:,}')
```

The total population is 8,644

```
In [62]:   1  # Calculate workforce total percentage
           2  workforce_total_popn = (workforce_age/total_population)*100
           3  print(f'The percentage of workforce age group is {workforce_total_popn:.2f}%')
```

The percentage of workforce age group is 66.01%

```
In [63]:   1  # Filter workforce_age from total unemployed
           2  employed_workforce_age = new_data[(new_data['Occupation Status'] == 'Employed')
           3                                    & (new_data['Age'] >= 18) & (new_data['Age'] <= 65)]['Occupation Status'].cou
           4
           5  # Filter workforce_age from total unemployed
           6  unemployed_workforce_age = new_data[(new_data['Occupation Status'] == 'Unemployed')
           7                                      & (new_data['Age'] >= 18) & (new_data['Age'] <= 65)]['Occupation Status'].cou
```

```python
In [64]:  1  # Check for employed that is within the workforce age
          2  employed_workforce = employed_workforce_age/(unemployed_workforce_age + employed_workforce_age)
          3
          4  # Multiply by 100 to get the rate
          5  employed_workforce_rate = employed_workforce * 100
          6
          7  # Print the result
          8  print(f'The employed workforce rate is {employed_workforce_rate:.2f}%')
```

The employed workforce rate is 89.43%

```python
In [65]:  1  # Check for unemployed that is within the workforce age
          2  unemployed_workforce = unemployed_workforce_age/(unemployed_workforce_age + employed_workforce_age)
          3
          4  # Multiply by 100 to get the rate
          5  unemployed_workforce_rate = unemployed_workforce * 100
          6
          7  # Print the result
          8  print(f'The unemployed workforce rate is {unemployed_workforce_rate:.2f}%')
```

The unemployed workforce rate is 10.57%

# Population Growth

## Birth Rate Calculation

```python
In [66]:
1  # Find the number of infant = 0 year old
2  infant = new_data['Age'] == 0
3  infant_num = len(new_data[infant])
4
5  # Filter women at age 25-29 years
6  age_25_29 = new_data[(new_data['Age'] >= 25) & (new_data['Age'] <= 29)]
7  female_age_25_29= len(age_25_29[age_25_29['Gender']=='Female'])
8
9  # Find the total birth of women between the age of 25_29 years
10 num_birth_25_29 = (infant_num/female_age_25_29)*100000
11 birth_rate_25_29 = (num_birth_25_29)/total_population
12
13 print(f'The birth rate /annum of women between 25 and 29 years is {birth_rate_25_29:.2f}%','of 100,000 population
```

The birth rate /annum of women between 25 and 29 years is 3.55% of 100,000 population:

```python
In [67]:
1  # Find the number of babies = 4 years old
2  babies = new_data['Age'] == 4
3  babies_num = len(new_data[babies])
4
5  # Filter women at age 30-34 years
6  age_30_34 = new_data[(new_data['Age'] >= 30) & (new_data['Age'] <= 34)]
7  fermale_age_30_34= len(age_30_34[age_30_34['Gender']=='Female'])
8
9  # Find the birth rate of women between the age of 30-34 years
10 num_birth_30_34 = (babies_num/fermale_age_30_34)*100000
11 birth_rate_30_34 = (num_birth_30_34)/total_population
12
13 print(f'The birth rate/annum of women between 30 and 34 years is {birth_rate_30_34:.2f}%','of 100,000 population:
```

The birth rate/annum of women between 30 and 34 years is 3.60% of 100,000 population:

## Death Rate Calculation

```python
In [68]:   1  # Filter the Age of persons between 50-59years
           2  aged_50_59 = new_data[(new_data['Age'] >= 50) & (new_data['Age'] <= 59)]
           3
           4  # Filter the Age of persons between 60-69years
           5  aged_60_69 = new_data[(new_data['Age'] >= 60) & (new_data['Age'] <= 69)]
           6
           7  # The difference between the two age groups
           8  death_1= len(aged_50_59) - len(aged_60_69)
           9
          10  # Base population = total population
          11  total_popn_1 =len(aged_50_59)
          12
          13  # Find the death rate per annum per 100,000 population
          14  no_of_death1 = (death_1/total_popn_1)*100000
          15
          16  # Number of death/annum
          17  annual_death_1 = int(no_of_death1/10)
          18
          19  # Print the result
          20  print(f'The death/annum of people betwwen 50 to 69 years is {annual_death_1:,}', 'of 100,000 population')
```

The death/annum of people betwwen 50 to 69 years is 3,517 of 100,000 population

```python
In [69]:  1  # Filter the Age of persons between 70-79years
          2  aged_70_79 = new_data[(new_data['Age'] >= 70) & (new_data['Age'] <= 79)]
          3
          4  # Filter the Age of persons between 80-89years
          5  aged_80_89 = new_data[(new_data['Age'] >= 80) & (new_data['Age'] <= 89)]
          6
          7  death_2 = len(aged_70_79)- len(aged_80_89)
          8  total_popn_2 = len(aged_70_79)
          9
         10  no_of_death2 = (death_2/total_popn_2)*100000
         11  annual_death_2 = int(no_of_death2/10)
         12
         13  # Print the result
         14  print(f'The death/annum of people betwwen 70 to 89 years is {annual_death_2:,}', 'of 100,000 population')
```

The death/annum of people betwwen 70 to 89 years is 5,811 of 100,000 population

```python
In [70]:  1  # Filter the Age of persons between 90-99years
          2  aged_90_99 = new_data[(new_data['Age'] >= 90) & (new_data['Age'] <= 99)]
          3
          4  # Filter the Age of persons between 100-110years
          5  aged_100_110 = new_data[(new_data['Age'] >= 100) & (new_data['Age'] <= 109)]
          6
          7  death_3 = len(aged_90_99) - len(aged_100_110)
          8  total_popn_3 = len(aged_90_99)
          9
         10  no_of_death3 = (death_3/total_popn_3)*100000
         11  annual_death_3 = int(no_of_death3/10)
         12
         13  # Print the result
         14  print(f'The death/annum of people betwwen 90 to 109 years is {annual_death_3:,}', 'of 100,000 population')
```

The death/annum of people betwwen 90 to 109 years is 833 of 100,000 population

```
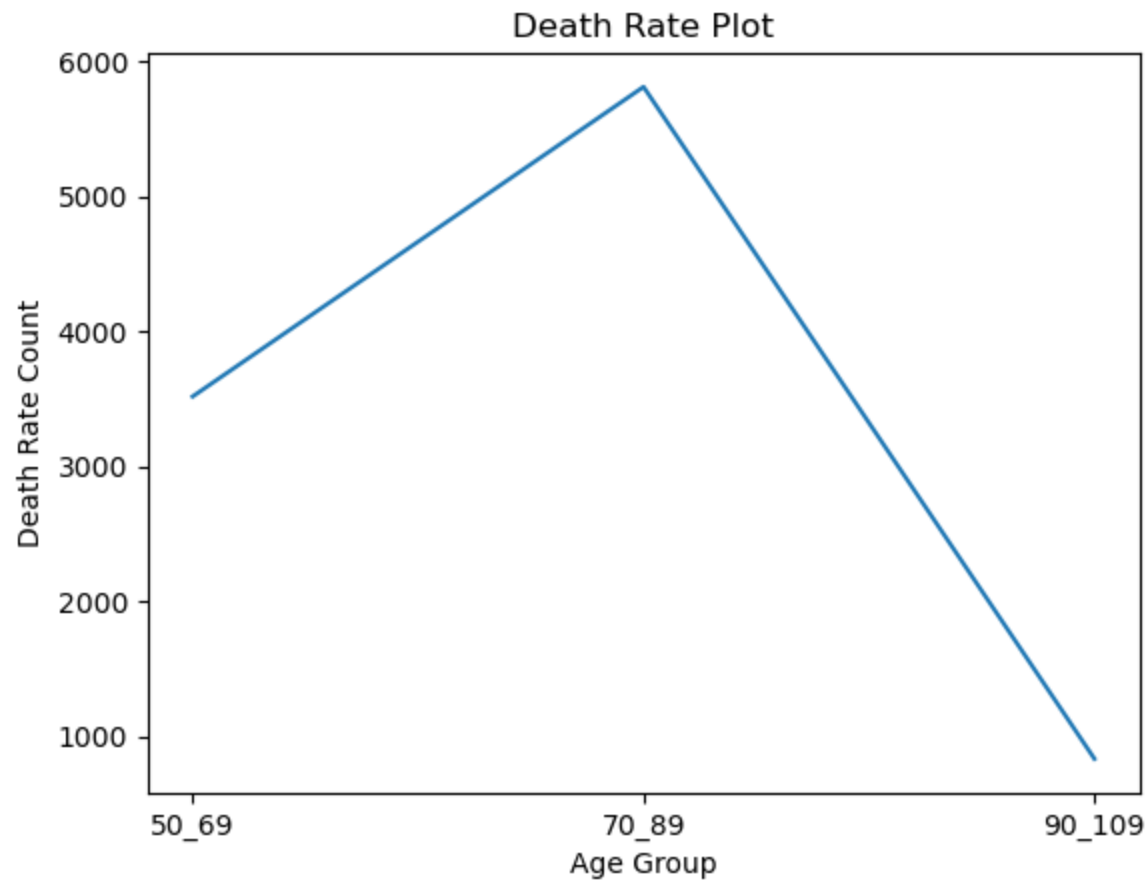In [71]:   1  # Total death of 100,000 population
           2  total_death = annual_death_1 + annual_death_2 + annual_death_3
           3
           4  # Total death rate of 100,000 population
           5  total_death_rate = total_death/total_population
           6  print(f'The total_death_rate/annum is {total_death_rate:.2f}%', 'of 100,000 population')
```

The total_death_rate/annum is 1.18% of 100,000 population

**Death Rate Plot**

```
In [72]:  1  # Line Plot for death rate
          2  death_rate_plot = sns.lineplot(x = ['50_69', '70_89', '90_109'], y = [annual_death_1, annual_death_2, annual_deat
          3
          4  # Title the plot, x and y axis
          5  death_rate_plot.set(title = 'Death Rate Plot')
          6  death_rate_plot.set_xlabel('Age Group')
          7  death_rate_plot.set_ylabel('Death Rate Count')
```

Out[72]:  Text(0, 0.5, 'Death Rate Count')

## Migration Calculation

```
In [73]:   1  # create a variable names for lodgers and visitors
           2  lodgers = new_data[new_data['Relationship to Head of House'] == 'Lodger']['Relationship to Head of House'].count(
           3  visitors = new_data[new_data['Relationship to Head of House'] == 'Visitor']['Relationship to Head of House'].coun
           4
           5  # create a variable names for uni_students and employed
           6  uni_students = new_data[new_data['Occupation'] == 'University Student']['Occupation'].count()
           7
           8  # create a variable name for divorced
           9  divorced = new_data[new_data['Marital Status'] == 'Divorced']['Marital Status'].count()
```

## Possible Immigrant

```
In [74]:   1  # Filter lodgers from total employed
           2  employed_lodgers = new_data[(new_data['Occupation Status'] == 'Employed')
           3                              & (new_data['Relationship to Head of House'] == 'Lodger')]['Occupation Status'].co
           4
           5  # Filter visitors from total employed
           6  employed_visitors = new_data[(new_data['Occupation Status'] == 'Employed')
           7                               & (new_data['Relationship to Head of House'] == 'Visitor')]['Occupation Status'].c
           8
           9  # Filter visitors from divorced
          10  divorced_visitors = new_data[(new_data['Marital Status'] == 'Divorced')
          11                               & (new_data['Relationship to Head of House'] == 'Visitor')]['Marital Status'].coun
          12
```

```
In [75]:   1  # Total number of possible immigrant
           2  immigrant = (employed_lodgers+ employed_visitors+ divorced_visitors)
           3  immigration_rate = immigrant/total_population
           4
           5  #Print the result
           6  print(f'The Immigration Rate/annum is {immigration_rate:.2f}%')
```

```
The Immigration Rate/annum is 0.04%
```

## Possible Emmigrant

```python
In [76]:  1  # Filter divorced from total unemployed
          2  unemployed_divorced = new_data[(new_data['Occupation Status'] == 'Unemployed')
          3                                 & (new_data['Marital Status'] == 'Divorced')]['Occupation Status'].count()
          4
          5  # Filter workforce_age from total unemployed
          6  unemployed_workforce_age = new_data[(new_data['Occupation Status'] == 'Unemployed')
          7                                      & (new_data['Age'] >= 18) & (new_data['Age'] <= 65)]['Occupation Status'].cou
          8
          9  # Filter divorced from lodgers
         10  divorced_lodgers = new_data[(new_data['Marital Status'] == 'Divorced')
         11                              & (new_data['Relationship to Head of House'] == 'Lodger')]['Marital Status'].count()
```

```python
In [77]:  1  #Total number of possible emmigrant
          2  emigrant = unemployed_divorced + unemployed_workforce_age + divorced_lodgers
          3  emigration_rate = emigrant/total_population
          4
          5  #Print the result
          6  print(f'The Emmigration Rate/annum is {emigration_rate:.2f}%')
```

The Emmigration Rate/annum is 0.08%

## Migration Rate

```python
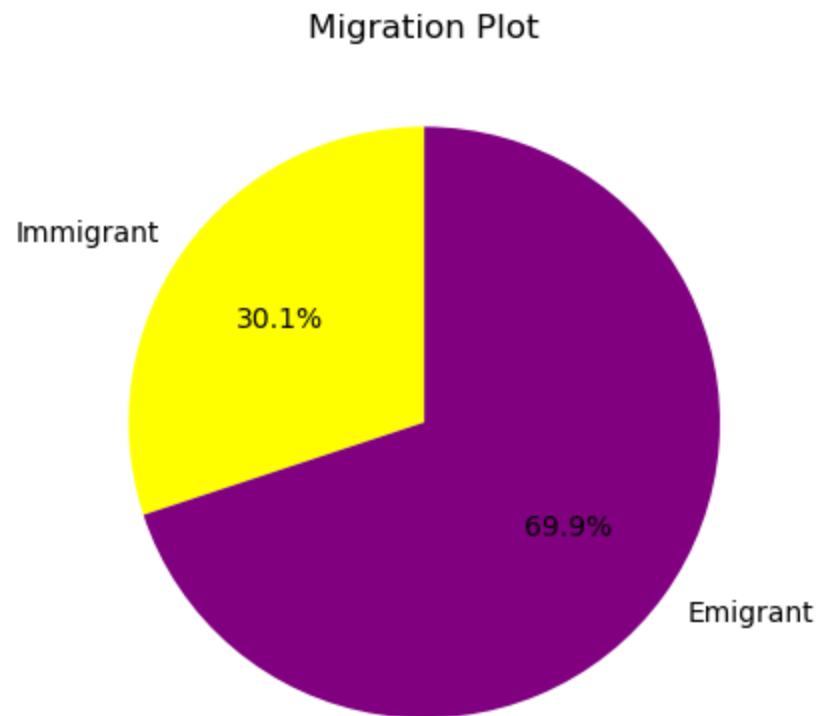In [78]:  1  # Calculate the migration rate
          2  migration_rate = immigration_rate - emigration_rate
          3
          4  # Print the result
          5  print(f'The migration_rate/annum is {migration_rate:.2f}%', 'of 100,000 apopulation')
```

The migration_rate/annum is -0.05% of 100,000 apopulation

```
In [79]:  1  # define the data
          2  labels = ['Immigrant', 'Emigrant']
          3  sizes = [immigration_rate, emigration_rate]
          4  colors = ['yellow','purple']
          5
          6  # create the pie chart
          7  plt.pie(sizes, labels=labels, colors=colors, autopct='%1.1f%%', startangle=90)
          8
          9  # title the plot
         10  plt.title('Migration Plot')
```

Out[79]:  Text(0.5, 1.0, 'Migration Plot')

## Population Growth Rate

```
In [80]:   1  # Define the Crude death and Crude birth rate
           2  crude_birth_rate = birth_rate_25_29
           3  crude_death_rate = total_death_rate
           4
           5  # Calculate population growth rate
           6  population_growth_rate = (crude_birth_rate + immigration_rate) - (crude_death_rate + emigration_rate)
           7
           8  # Print the result
           9  print(f'population_growth_rate/annum is {population_growth_rate:.2f}%', 'of 100,000 population')
```

population_growth_rate/annum is 2.32% of 100,000 population

# Occupancy Level

```
In [81]:  1  # Find the number of houses, grouped by 'House Number and Street
          2  num_of_houses = new_data.groupby(['House Number', 'Street']).size().reset_index(name = 'Actual Occupant')
          3
          4  # Exclude house number 1 (It is assumed to be an hotel)
          5  num_of_houses = num_of_houses[~(num_of_houses['House Number'] == 1)]
          6  num_of_houses
```

Out[81]:

| | House Number | Street | Actual Occupant |
|---|---|---|---|
| **105** | 2 | Albans Pines | 5 |
| **106** | 2 | Albionfold Road | 4 |
| **107** | 2 | Archer Drives | 6 |
| **108** | 2 | Arrows Camp | 5 |
| **109** | 2 | Atkinson Meadows | 3 |
| **...** | ... | ... | ... |
| **2900** | 242 | Hughes Camp | 2 |
| **2901** | 243 | Hughes Camp | 2 |
| **2902** | 244 | Hughes Camp | 2 |
| **2903** | 245 | Hughes Camp | 2 |
| **2904** | 246 | Hughes Camp | 2 |

```
In [82]:  1  # Output the total number of houses
          2  total_houses = num_of_houses.shape[0]
          3
          4  print(f"The number of houses are {total_houses:.0f}")
```

The number of houses are 2800

```
In [83]:   1   # Find the average occupant staying in each house
           2   average_occupant = num_of_houses['Actual Occupant'].mean()
           3
           4   print(f"The Average Occupant is {average_occupant:.0f}")
```

The Average Occupant is 3

```
In [84]:   1   # Number of houses whose occupant are less than or equal to the average occupant
           2   occupant_less_than_3 = num_of_houses[num_of_houses['Actual Occupant']<= 3].shape[0]
           3
           4   print(f"The Occupant less than or equal to 3 is {occupant_less_than_3:.0f}")
```

The Occupant less than or equal to 3 is 1806

```
In [85]:   1   # Number of houses whose occupant are greater than average occupant
           2   occupant_greater_than_3 = num_of_houses[num_of_houses['Actual Occupant']> 3].shape[0]
           3
           4   print(f"The Occupant greater than 3 is {occupant_greater_than_3:.0f}")
```

The Occupant greater than 3 is 994

```
In [86]:   1   # Percentage of houses greater than average occupant
           2   houses_overused = (occupant_greater_than_3/total_houses)*100
           3
           4   print(f"The percentage of houses overused are: {houses_overused:.0f}%")
```

The percentage of houses overused are: 36%

## Possible Commuters

```
In [87]:   1   phd_students = new_data['Occupation'].str.contains('phd', case=False)
           2   phd_students = phd_students.sum()
```

```
In [88]:   1  lecturers = new_data['Occupation'].str.contains('lecturer', case=False) & ~new_data['Occupation'].str.contains('r
           2  lecturers = lecturers.sum()
```

```
In [89]:   1  professors = new_data['Occupation'].str.contains('professor', case=False) & ~new_data['Occupation'].str.contains(
           2  professors = professors.sum()
```

```
In [90]:   1  academic_staff = new_data['Occupation'].str.contains('academic', case=False) & ~new_data['Occupation'].str.contai
           2  academic_staff = academic_staff.sum()
```

```
In [91]:   1  # Total number of commuters
           2  commuters = uni_students + phd_students + lecturers + professors + academic_staff
           3
           4  print(f'The total number of commuters is {commuters:.0f}')
```

The total number of commuters is 661

```
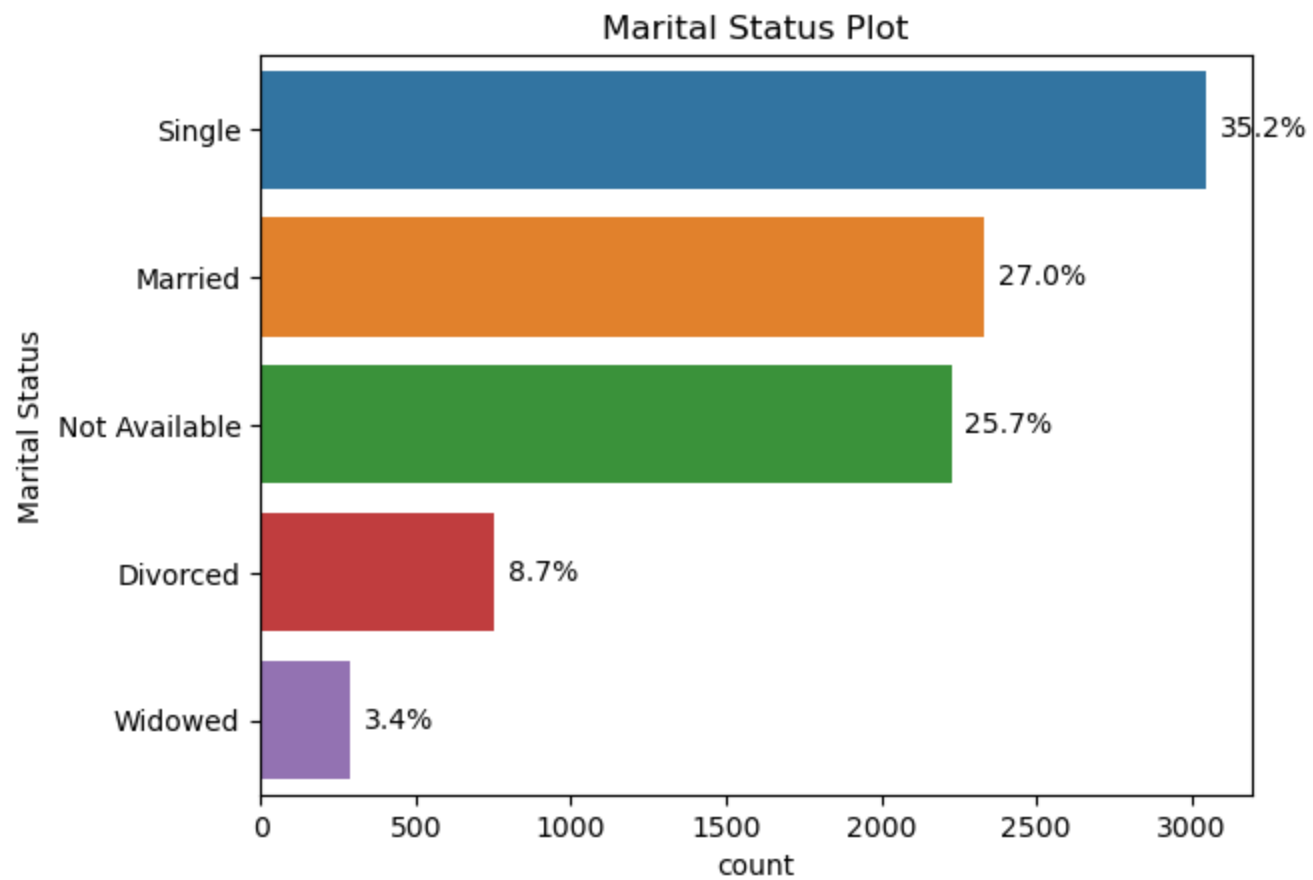In [92]:   1  # Percentage of commuters
           2  perc_commuters = (commuters/total_population) *100
           3
           4  # Print the result
           5  print(f'The percentage of commuters is {perc_commuters:.2f}%')
```

The percentage of commuters is 7.65%

## Marital Status

```
In [93]:    1  # Create a plot for the Gender column
            2  marital_plot = sns.countplot(data=new_data, y='Marital Status')
            3
            4  # Get the total frequency of each category
            5  total_frequency = len(new_data['Marital Status'])
            6
            7  # Loop through each bar in the plot
            8  for p in marital_plot.patches:
            9      # Get the percentage value for the bar
           10      percentage = p.get_width() / total_frequency * 100
           11
           12      # Add the percentage as text next to the bar
           13      marital_plot.annotate(f'{percentage:.1f}%', (p.get_width(), p.get_y() + p.get_height() / 2.), xytext=(5, 0),
           14
           15  # Rotate the y-axis label
           16  marital_plot.set_yticklabels(marital_plot.get_yticklabels(), rotation=0)
           17
           18  # Give the plot a title and name the y-axis
           19  marital_plot.set(title='Marital Status Plot')
           20  marital_plot.set_ylabel('Marital Status')
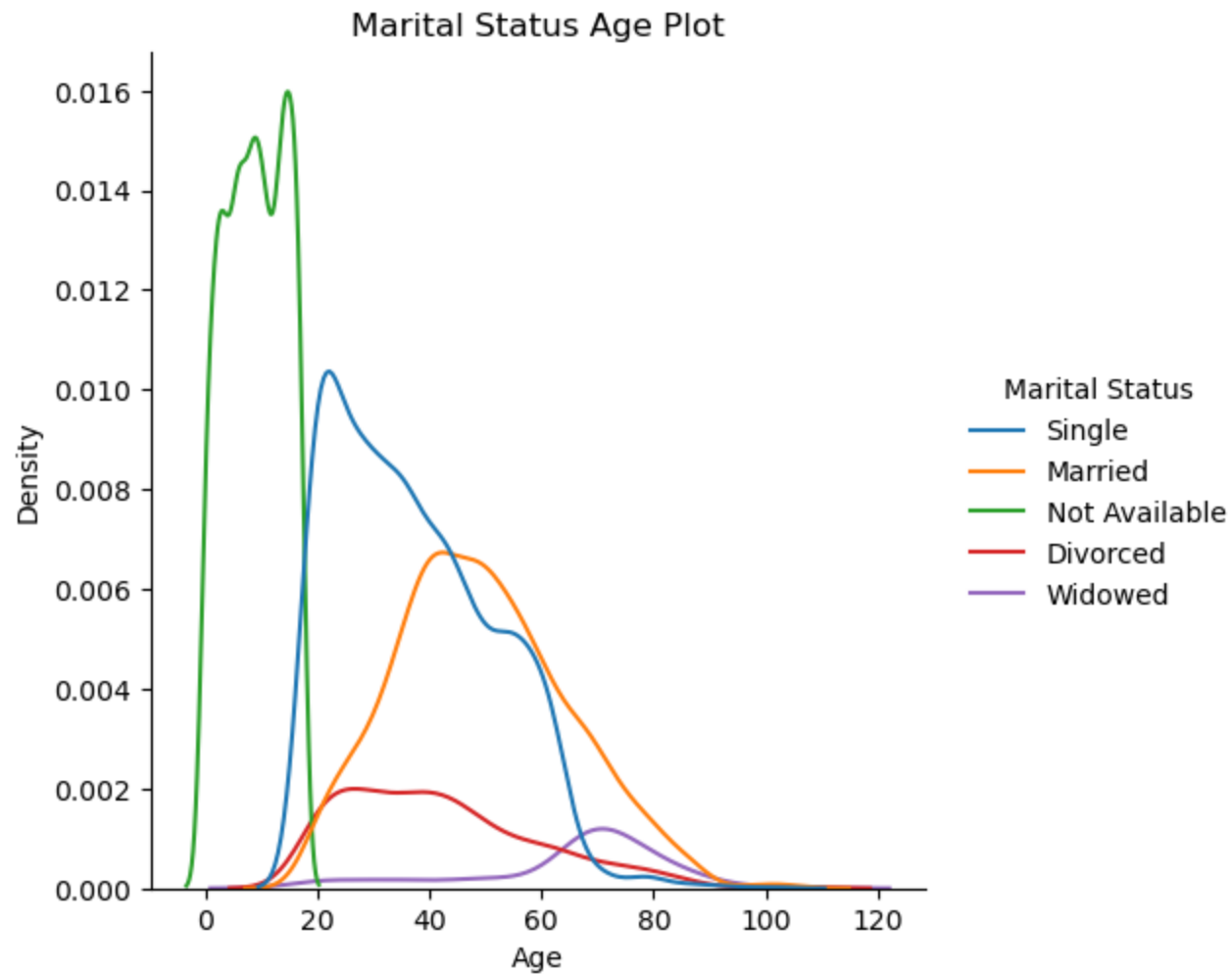```

Out[93]:   Text(0, 0.5, 'Marital Status')

Marital Status Plot

```python
# Create a plot for the Gender column
marital_plot = sns.countplot(data=new_data, y='Marital Status', hue='Gender')

# Get the total frequency of each category
total_frequency = len(new_data['Marital Status'])

# Loop through each bar in the plot
for p in marital_plot.patches:
    # Get the percentage value for the bar
    percentage = p.get_width() / total_frequency * 100

    # Add the percentage as text next to the bar
    marital_plot.annotate(f'{percentage:.1f}%', (p.get_width(), p.get_y() + p.get_height() / 2.), xytext=(5, 0),

# Rotate the y-axis label
marital_plot.set_yticklabels(marital_plot.get_yticklabels(), rotation=0)

# Give the plot a title and name the y-axis
marital_plot.set(title='Marital Status_Gender Plot')
marital_plot.set_ylabel('Marital Status')
```

Out[94]: Text(0, 0.5, 'Marital Status')

```
1  # Plot for age and marital status
2  marital_age_plot = sns.displot(new_data, x="Age", hue="Marital Status", kind="kde")
3
4  # Title the plot
5  marital_age_plot.set(title = 'Marital Status Age Plot')
```

Out[95]: &lt;seaborn.axisgrid.FacetGrid at 0x1cab7b13ee0&gt;

```
In [96]:  1  # Calculate divorce rate against the total population
          2  divorced_rate = (divorced/total_population)*100
          3
          4  # Print the result
          5  print(f'The divorce rate is {divorced_rate:.2f}%')
```

The divorce rate is 8.72%

```
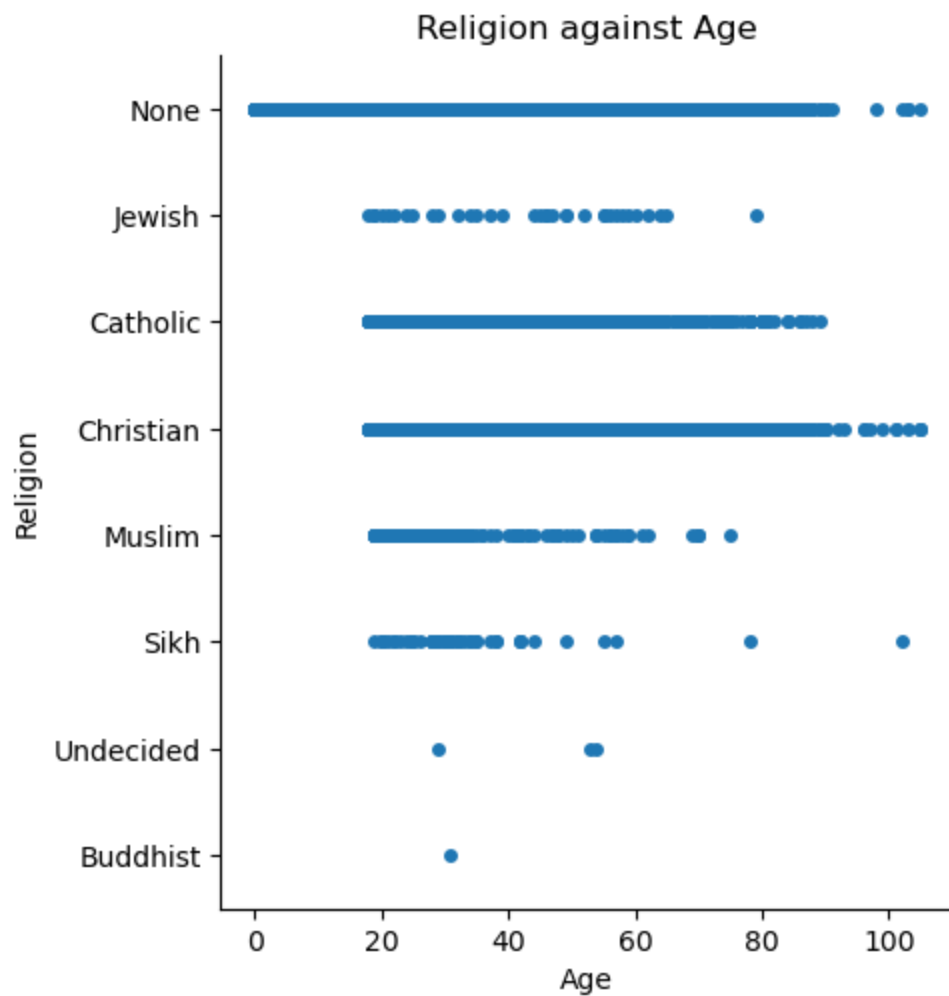In [97]:  1  # create a variable for married people
          2  married = new_data[new_data['Marital Status'] == 'Married']['Marital Status'].count()
          3
          4  # Calculate divorce rate against the total population
          5  married_rate = (married/total_population)*100
          6
          7  # Print the resultbirth
          8  print(f'The married rate is {married_rate:.2f}%')
```

The married rate is 26.96%

## Religious Affliations

```python
# Plot for religion with age
religion_catplot = sns.catplot(new_data, x="Age", y="Religion", jitter=False)

# Title the plot
religion_catplot.set(title = 'Religion against Age')
```

Out[98]: <seaborn.axisgrid.FacetGrid at 0x1cab909d400>

Religion against Age

```
In [99]:  1  # Create a plot for the religion
          2  religion_countplot = sns.countplot(data=new_data, y='Religion')
          3
          4  # Get the total frequency of each category
          5  total_frequency = len(new_data['Religion'])
          6
          7  # Loop through each bar in the plot
          8  for p in religion_countplot.patches:
          9      # Get the percentage value for the bar
         10      percentage = p.get_width() / total_frequency * 100
         11
         12      # Add the percentage as text next to the bar
         13      religion_countplot.annotate(f'{percentage:.1f}%', (p.get_width(), p.get_y() + p.get_height() / 2.), xytext=(5
         14
         15  # Give the plot a title and name the y-axis
         16  religion_countplot.set(title='Relion Affliation Plot')
```

Out[99]: [Text(0.5, 1.0, 'Relion Affliation Plot')]