# Task 1: Data Preparation (Primary.csv only, the others are prepared for tasks)

```
In [1]:  1  import pandas as pd
         2  primary_filepathog = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_Data
         3  pogdf = pd.read_csv(primary_filepathog, encoding='latin1')
```

```
In [2]:  1  pogdf.head()
```

Out[2]:

| | column A | column B | column C | column D | column E | column F | column G | column H | column I | column J | column K | colu |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | T pe |
| 1 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 15% | 2% | 22% | 0% | 61% | Demographic and Health Survey | 20 |
| 2 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 39% | NaN | NaN | NaN | NaN | Multiple Indicator Cluster Survey | 20 |
| 3 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 81% | 69% | 89% | 46% | 99% | Demographic and Health Survey | 20 |
| 4 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 34% | 30% | 49% | 7% | 75% | Multiple Indicator Cluster Survey | 2 |

```
In [3]:  1  pogdf.isnull().any()
```

```
Out[3]:  column A     False
         column B     False
         column C     False
         column  D    False
         column E     False
         column F     False
         column G      True
         column H      True
         column I      True
         column J      True
         column K     False
         column L     False
         dtype: bool
```

```
In [4]:  1  duplicates = pogdf.duplicated()
         2  print(f"Number of duplicate rows: {duplicates.sum()}")
```

```
Number of duplicate rows: 0
```

```
In [5]:  1  #now that we know which columns have missing values, I want to check how many missing
         2  missing_values = pogdf.isnull().sum()
         3  print("Missing values per column:\n", missing_values)
```

```
Missing values per column:
 column A      0
column B      0
column C      0
column  D     0
column E      0
column F      0
column G     11
column H      8
column I     18
column J     21
column K      0
column L      0
dtype: int64
```

```
In [6]:   1  #here is how I will be cleaning the data
          2  #first I will pull the csv into a dataframe
          3  #i till skip the first row as it is just the columns by letters which isn't the most
          4  #to remember or analyze
          5
          6  file_path_test = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSciV
          7  dftest = pd.read_csv(file_path_test, skiprows=1)
          8
          9  output_file_path_test = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_I
         10  dftest.to_csv(output_file_path_test, index=False)
         11  dftest.head()
```

Out[6]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Time period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 15% | 2% | 22% | 0% | 61% | Demographic and Health Survey | 2015-16 |
| 1 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 39% | NaN | NaN | NaN | NaN | Multiple Indicator Cluster Survey | 2011-12 |
| 2 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 81% | 69% | 89% | 46% | 99% | Demographic and Health Survey | 2015-16 |
| 3 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 34% | 30% | 49% | 7% | 75% | Multiple Indicator Cluster Survey | 2019 |
| 4 | BRB | Barbados | LAC | LAC | High income (H) | 63% | 54% | 68% | 9% | 97% | Multiple Indicator Cluster Survey | 2012 |

```python
import numpy as np

file_path_test = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSciV
dftest = pd.read_csv(file_path_test, skiprows=1)

# Function for converting the percentage string value to float
def percentage_to_float(x):
    if isinstance(x, str) and x.endswith('%'):
        return float(x[:-1]) / 100
    return x

# Converting the string percentage values to float
cols_to_convert = ['Total',
                   'Rural (Residence)',
                   'Urban (Residence)',
                   'Poorest (Wealth quintile)',
                   'Richest (Wealth quintile)']
for col in cols_to_convert:
    dftest[col] = dftest[col].apply(percentage_to_float)

# Replacing the empty 'NaN' values in 'Rural (Residence)' and 'Urban (Residence)'
#columns with the mean of other non-empty countries in the same 'Region'
cols_to_fill = ['Rural (Residence)',
                'Urban (Residence)',
                'Poorest (Wealth quintile)',
                'Richest (Wealth quintile)']

for col in cols_to_fill:
    dftest[col] = dftest.groupby('Region')[col].apply(lambda x: x.fillna(round(x.mean

output_file_path_test = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_[
dftest.to_csv(output_file_path_test, index=False)
dftest.head()
```

Out[7]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Time period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.15 | 0.02 | 0.22 | 0.00 | 0.61 | Demographic and Health Survey | 2015-16 |
| 1 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.39 | 0.26 | 0.47 | 0.22 | 0.80 | Multiple Indicator Cluster Survey | 2011-12 |
| 2 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 | 0.69 | 0.89 | 0.46 | 0.99 | Demographic and Health Survey | 2015-16 |
| 3 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.34 | 0.30 | 0.49 | 0.07 | 0.75 | Multiple Indicator Cluster Survey | 2019 |
| 4 | BRB | Barbados | LAC | LAC | High income (H) | 0.63 | 0.54 | 0.68 | 0.09 | 0.97 | Multiple Indicator Cluster Survey | 2012 |

```python
#using regex to handle the 'Time period' values in the last column
#basically to ensure all he values in the column are in the 2010's
import re

def handle_time_period(value):
    # Case 1: If format is like '2015-16' or '2011-12'
    if re.match(r'\d{4}-\d{2}', value):
        return value[:2] + value[-2:]

    # Case 2: If the value is something outlandish like '2562'
    if int(value) > 2019 or int(value) < 2010:
        return '20' + value[-2:]

    # Case 3: If the value is like '2027'
    if int(value) > 2019:
        return '201' + value[-1]

    # Case 4: If the value is like '2012-99'
    if re.match(r'\d{4}-\d{2}', value) and int(value[-2:]) > 19:
        return '20' + value[-2:]

    # Case 5: If the value is like '2018-2019'
    if re.match(r'\d{4}-\d{4}', value):
        return value[-4:]

    # Case 6: If the value is like '2076'
    if int(value) > 2019 and int(value[-1]) > 0:
        return '20' + value[-1]

    return value

# Applying the function to the 'Time period' column
dftest['Time period'] = dftest['Time period'].apply(handle_time_period)

output_file_path_test = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_I
dftest.to_csv(output_file_path_test, index=False)

dftest.head(59)
```

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **29** | GNB | Guinea-Bissau | SSA | WCA | Low income (L) | 0.02 | 0.01 | 0.04 | 0.00 | 0.05 | Multiple Indicator Cluster Survey | 2019 |
| **30** | HTI | Haiti | LAC | LAC | Low income (L) | 0.18 | 0.10 | 0.33 | 0.00 | 0.60 | Demographic and Health Survey | 2017 |
| **31** | IND | India | SA | SA | Lower middle income (LM) | 0.07 | 0.04 | 0.14 | 0.00 | 0.33 | Demographic and Health Survey | 2016 |
| **32** | IDN | Indonesia | EAP | EAP | Lower middle income (LM) | 0.17 | 0.09 | 0.24 | 0.05 | 0.44 | SUSENAS | 2019 |
| **33** | IRQ | Iraq | MENA | MENA | Upper middle income (UM) | 0.46 | 0.33 | 0.53 | 0.14 | 0.83 | Multiple Indicator Cluster Survey | 2018 |

```
In [9]:  1  #I am removing the following columns as I feel they aren't relevant to the analysis
         2  columns_to_remove1 = ['Rural (Residence)','Urban (Residence)','Poorest (Wealth quint:
         3  dftest = dftest.drop(columns=columns_to_remove1)
         4  dftest.head()
```

Out[9]:

|   | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|---|------|---------------------|--------|------------|--------------|-------|
| 0 | AGO  | Angola              | SSA    | ESA        | Lower middle income (LM) | 0.15 |
| 1 | ARG  | Argentina           | LAC    | LAC        | Upper middle income (UM) | 0.39 |
| 2 | ARM  | Armenia             | ECA    | EECA       | Upper middle income (UM) | 0.81 |
| 3 | BGD  | Bangladesh          | SA     | SA         | Lower middle income (LM) | 0.34 |
| 4 | BRB  | Barbados            | LAC    | LAC        | High income (H)          | 0.63 |

```
In [10]:  1  #now that I have the dataset I wanted, i will convert the string values into floats
          2  #this will allow me to analyze the data numerically
          3
          4  def percentage_to_float(value):
          5      if isinstance(value, str) and value.endswith('%'):
          6          return float(value[:-1]) / 100
          7      else:
          8          return value
          9
         10  columns_to_convert = ['Total']
         11
         12  for column in columns_to_convert:
         13      dftest[column] = dftest[column].apply(percentage_to_float)
         14
         15  dftest.head()
```

Out[10]:

|   | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|---|------|---------------------|--------|------------|--------------|-------|
| 0 | AGO  | Angola              | SSA    | ESA        | Lower middle income (LM) | 0.15 |
| 1 | ARG  | Argentina           | LAC    | LAC        | Upper middle income (UM) | 0.39 |
| 2 | ARM  | Armenia             | ECA    | EECA       | Upper middle income (UM) | 0.81 |
| 3 | BGD  | Bangladesh          | SA     | SA         | Lower middle income (LM) | 0.34 |
| 4 | BRB  | Barbados            | LAC    | LAC        | High income (H)          | 0.63 |

## Task 2.1: I am choosing the following:

Nominal Value: 'Region; Ordinal Value: 'Income Group'; Numerical Value: 'Total';

```
In [11]:  1  #my data is now cleaned as I have the rows that I want to analyze
          2  import seaborn as sns
          3  import matplotlib.pyplot as plt
```

```
In [12]:   1  plt.figure(figsize=(12, 6))
           2  sns.boxplot(x='Region', y='Total', data=dftest)
           3  plt.title('Box Plot of Total by Region')
           4  plt.xlabel('Region')
           5  plt.ylabel('Total')
           6  plt.show()
```



Box Plot of Total by Region

```
In [13]:   1  #As we can witness, there are outliers on the 'SSA' region
           2
           3  def detect_outliers(region_df):
           4      Q1 = region_df['Total'].quantile(0.25)
           5      Q3 = region_df['Total'].quantile(0.75)
           6      IQR = Q3 - Q1
           7
           8      lower_bound = Q1 - 1.5 * IQR
           9      upper_bound = Q3 + 1.5 * IQR
          10
          11      outliers = region_df[(region_df['Total'] < lower_bound) | (region_df['Total'] > u
          12      return outliers
          13
          14  region_groups = dftest.groupby('Region')
          15
          16  outliers = pd.concat([detect_outliers(group) for _, group in region_groups])
          17
          18  outliers.head()
```

Out[13]:

|    | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|----|------|---------------------|--------|------------|--------------|-------|
| 25 | GMB  | Gambia              | SSA    | WCA        | Low income (L) | 0.64 |
| 60 | STP  | Sao Tome and Principe | SSA  | WCA        | Lower middle income (LM) | 0.44 |

```
#I want to explore and analyze these so we can make a good inference as to what make
#To do this, I am going to visualize a similar boxplot, with the Sub-region and Tota
#my reasoning for this is that maybe Sub-region migh be a better predictor of the To
# and Income Group

plt.figure(figsize=(12, 6))
sns.boxplot(x='Sub-region', y='Total', data=dftest)
plt.title('Box Plot of Total by Sub-region')
plt.xlabel('Sub-region')
plt.ylabel('Total')
plt.show()
```



Box Plot of Total by Sub-region

```
In [15]:    1  plt.figure(figsize=(12, 6))
            2  sns.boxplot(x='Region', y='Total', data=dftest)
            3  plt.title('Box Plot of Total by Region')
            4  plt.xlabel('Region')
            5  plt.ylabel('Total')
            6  plt.show()
            7
            8  plt.figure(figsize=(12, 6))
            9  sns.boxplot(x='Sub-region', y='Total', data=dftest)
           10  plt.title('Box Plot of Total by Sub-region')
           11  plt.xlabel('Sub-region')
           12  plt.ylabel('Total')
           13  plt.show()
           14
           15  plt.figure(figsize=(12, 6))
           16  sns.boxplot(x='Income Group', y='Total', data=dftest)
           17  plt.title('Box Plot of Total by Income Group')
           18  plt.xlabel('Income Group')
           19  plt.ylabel('Total')
           20  plt.show()
```



Box Plot of Total by Region

Box Plot of Total by Sub-region



Box Plot of Total by Income Group

In [16]:
```
1  #I wanted to see he spread of data across the Primary school children data
2  #for this, I employed a Histogram
3  #as we can see, most of the data is aggregated towards only 1% of of primary schoolc
4  #internet access
5
6
7  plt.figure(figsize=(12, 6))
8  sns.histplot(data=dftest, x='Total', bins=10)
9  plt.title('Histogram of Total Internet Access by Primary School Children')
10 plt.xlabel('Total')
11 plt.ylabel('Frequency')
12 plt.show()
13
```



Histogram of Total Internet Access by Primary School Children

```
In [17]:   1  plt.figure(figsize=(12, 6))
           2  ax1 = sns.barplot(x='Region', y='Total', data=dftest)
           3  plt.title('Bar Plot of Mean Total Internet Access by Region')
           4  plt.xlabel('Region')
           5  plt.ylabel('Mean Total')
           6
           7  # Setting the y-axis range from 0 to 1 to accomodate percent values
           8  ax1.set_ylim(0, 1)
           9
          10  plt.show()
          11
          12  plt.figure(figsize=(12, 6))
          13  ax2 = sns.barplot(x='Sub-region', y='Total', data=dftest)
          14  plt.title('Bar Plot of Mean Total Internet Access by Sub-region')
          15  plt.xlabel('Sub-region')
          16  plt.ylabel('Mean Total')
          17
          18  ax2.set_ylim(0, 1)
          19
          20  plt.show()
          21
          22  plt.figure(figsize=(12, 6))
          23  ax3 = sns.barplot(x='Income Group', y='Total', data=dftest)
          24  plt.title('Bar Plot of Mean Total Internet Access by Income Group')
          25  plt.xlabel('Income Group')
          26  plt.ylabel('Mean Total')
          27
          28  ax3.set_ylim(0, 1)
          29
          30  plt.show()
          31
```



Bar Plot of Mean Total Internet Access by Region

Bar Plot of Mean Total Internet Access by Sub-region

Bar Plot of Mean Total Internet Access by Income Group

```
In [18]:   1  plt.figure(figsize=(12, 6))
           2  ax1 = sns.violinplot(x='Region', y='Total', data=dftest)
           3  plt.title('Violin Plot of Total Internet Access by Region')
           4  plt.xlabel('Region')
           5  plt.ylabel('Total')
           6
           7  # Settting the y-axis range from -0.30 to 1.50 to accomodate the full violin plot
           8  ax1.set_ylim(-0.40, 1.50)
           9
          10  plt.show()
          11
          12  plt.figure(figsize=(12, 6))
          13  ax2 = sns.violinplot(x='Sub-region', y='Total', data=dftest)
          14  plt.title('Violin Plot of Total Internet Access by Sub-region')
          15  plt.xlabel('Region')
          16  plt.ylabel('Total')
          17
          18  ax2.set_ylim(-0.40, 1.50)
          19
          20  plt.show()
          21
          22  plt.figure(figsize=(12, 6))
          23  ax3 = sns.violinplot(x='Income Group', y='Total', data=dftest)
          24  plt.title('Violin Plot of Total Internet Access by Income Group')
          25  plt.xlabel('Income Group')
          26  plt.ylabel('Total')
          27
          28  ax3.set_ylim(-0.40, 1.50)
          29
          30  plt.show()
          31
```



Violin Plot of Total Internet Access by Region

Violin Plot of Total Internet Access by Sub-region

Violin Plot of Total Internet Access by Income Group

```
In [19]:   1  plt.figure(figsize=(12, 6))
           2  sns.swarmplot(x='Region', y='Total', data=dftest)
           3  plt.title('Swarm Plot of Total Internet Access by Region')
           4  plt.xlabel('Region')
           5  plt.ylabel('Total')
           6  plt.show()
           7
           8  plt.figure(figsize=(12, 6))
           9  sns.swarmplot(x='Sub-region', y='Total', data=dftest)
          10  plt.title('Swarm Plot of Total Internet Access by Sub-region')
          11  plt.xlabel('Sub-region')
          12  plt.ylabel('Total')
          13  plt.show()
          14
          15  plt.figure(figsize=(12, 6))
          16  sns.swarmplot(x='Income Group', y='Total', data=dftest)
          17  plt.title('Swarm Plot of Total Internet Access by Income Group')
          18  plt.xlabel('Income Group')
          19  plt.ylabel('Total')
          20  plt.show()
```
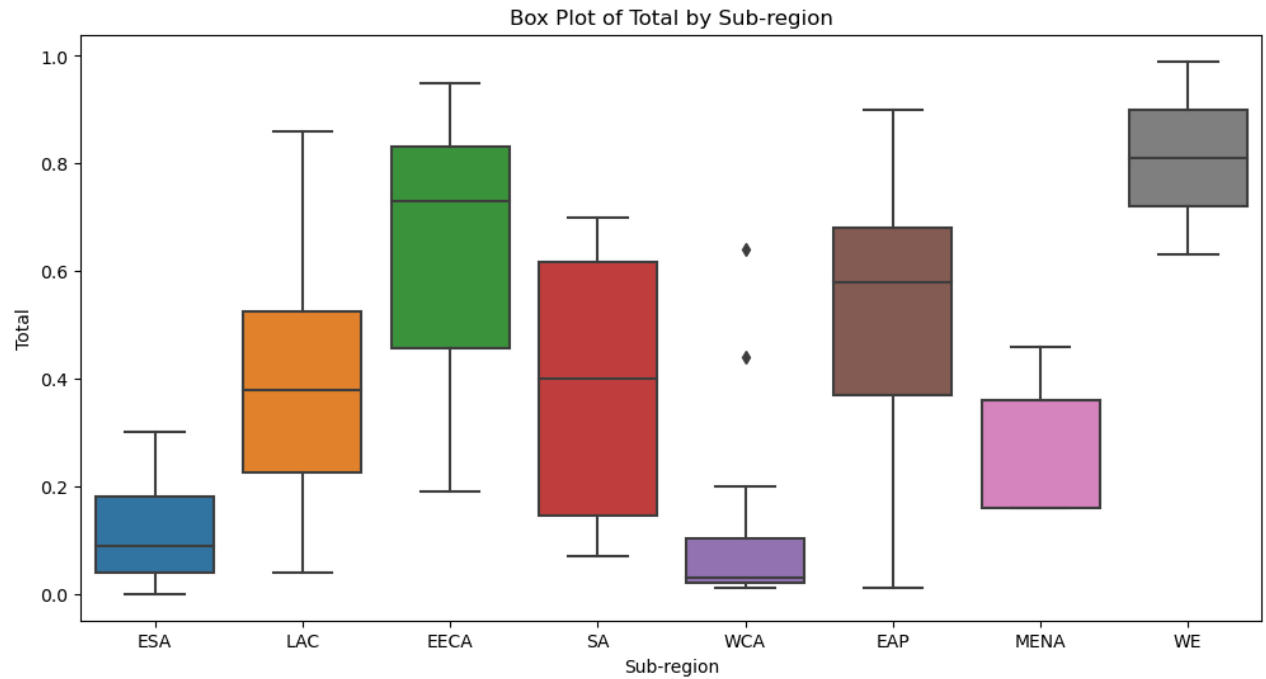


Swarm Plot of Total Internet Access by Region

Swarm Plot of Total Internet Access by Sub-region

Swarm Plot of Total Internet Access by Income Group

```
In [20]:    1  #I am using a Jittered Plot to see a messier version of the swarm plot so i can get
            2  #an idea of how much variation there is in between the data
            3
            4
            5  plt.figure(figsize=(12, 6))
            6  sns.stripplot(x='Region', y='Total', data=dftest, jitter=True, edgecolor='gray')
            7  plt.title('Jittered Scatter Plot of Total by Region')
            8  plt.xlabel('Region')
            9  plt.ylabel('Total')
           10  plt.show()
           11
           12  plt.figure(figsize=(12, 6))
           13  sns.stripplot(x='Sub-region', y='Total', data=dftest, jitter=True, edgecolor='gray')
           14  plt.title('Jittered Scatter Plot of Total by Sub-region')
           15  plt.xlabel('Sub-region')
           16  plt.ylabel('Total')
           17  plt.show()
           18
           19  plt.figure(figsize=(12, 6))
           20  sns.stripplot(x='Income Group', y='Total', data=dftest, jitter=True, edgecolor='gray
           21  plt.title('Jittered Scatter Plot of Total by Income Group')
           22  plt.xlabel('Income Group')
           23  plt.ylabel('Total')
           24  plt.show()
           25
```



Jittered Scatter Plot of Total by Region

## Task 1: ANOVA ANALYSIS FOR NUMERICAL VALUE SUBSTITUTION

Here, I am evaluating the effectiveness of 'Total', 'Region', and 'Income Group' by virtue of their 'Total' value predictive power

```python
In [21]:   1  #anova analysis to determine which is a better predictor of 'Total', 'Region', 'Sub-r
           2
           3  from scipy.stats import f_oneway
           4
           5  # ANOVA calculation for 'Region' and 'Total'
           6  region_groups = dftest.groupby('Region')['Total'].apply(list)
           7  region_anova_results = f_oneway(*region_groups)
           8  print(f"Region - F statistic: {region_anova_results.statistic:.2f}, P-value: {region_
           9
          10  #'Sub-region' and 'Total'
          11  sub_region_groups = dftest.groupby('Sub-region')['Total'].apply(list)
          12  sub_region_anova_results = f_oneway(*sub_region_groups)
          13  print(f"Sub-region - F statistic: {sub_region_anova_results.statistic:.2f}, P-value:
          14
          15  #'Income Group' and 'Total'
          16  income_group_groups = dftest.groupby('Income Group')['Total'].apply(list)
          17  income_group_anova_results = f_oneway(*income_group_groups)
          18  print(f"Income Group - F statistic: {income_group_anova_results.statistic:.2f}, P-val
          19
```

```
Region - F statistic: 15.30, P-value: 0.00000
Sub-region - F statistic: 10.93, P-value: 0.00000
Income Group - F statistic: 20.17, P-value: 0.00000
```

```python
In [22]:   1  #eta squared test to determine prediction and effective power of a categorical varial
           2  def eta_squared(f_statistic, df_between, df_within):
           3      return f_statistic * df_between / (f_statistic * df_between + df_within)
           4
           5  region_eta_squared = eta_squared(region_anova_results.statistic, len(region_groups)
           6  income_group_eta_squared = eta_squared(income_group_anova_results.statistic, len(inco
           7  sub_region_eta_squared = eta_squared(sub_region_anova_results.statistic, len(sub_reg:
           8
           9
          10  print(f"Region - Eta-squared: {region_eta_squared:.4f}")
          11  print(f"Income Group - Eta-squared: {income_group_eta_squared:.4f}")
          12  print(f"Sub-region - Eta-squared: {sub_region_eta_squared:.4f}")
          13
```

```
Region - Eta-squared: 0.4857
Income Group - Eta-squared: 0.4216
Sub-region - Eta-squared: 0.4921
```

```
In [23]:    1  #from the above Eta-squared value, We have inferred that the Sub-region is the best
            2  #'total' percent of children, however, I want to undersand how feasible this is to be
            3  #as a substitute for the mean if one seems o be missing, given this need, I feel I w
            4  #if in case, I want to compute and substitute a value for a missing value in a regio
            5
            6  region_counts = dftest['Region'].value_counts()
            7  print("Region Counts:")
            8  print(region_counts)
            9
           10  subregion_counts = dftest['Sub-region'].value_counts()
           11  print("\nSub-region Counts:")
           12  print(subregion_counts)
           13
           14  income_group_counts = dftest['Income Group'].value_counts()
           15  print("\nIncome Group Counts:")
           16  print(income_group_counts)
           17
```

```
Region Counts:
SSA      31
LAC      20
ECA      16
EAP       9
SA        6
MENA      5
Name: Region, dtype: int64

Sub-region Counts:
LAC      20
WCA      18
EECA     14
ESA      13
EAP       9
SA        6
MENA      5
WE        2
Name: Sub-region, dtype: int64

Income Group Counts:
Upper middle income (UM)      31
Lower middle income (LM)      30
Low income (L)                18
High income (H)                8
Name: Income Group, dtype: int64
```

## Task 2.2: Total School Age.csv data preparation and visualization

In this section, I will clean and substitute the missing values in the Total School Age csv I will employ the same methods that I used for preparing Primary.csv in Task 1

```
In [24]:  signment1_DataSciWithPython_s3970066/Global database on school-age digital connectivity-C
```

```
In [25]:    1  df3 = pd.read_csv(file_path3, encoding='latin1')
```

```
In [26]:    1  df3.describe()
```

Out[26]:

| | column A | column B | column C | column D | column E | column F | column G | column H | column I | column J | column K | column L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 88 | 88 | 88 | 88 | 88 | 88 | 78 | 81 | 71 | 70 | 88 | 88 |
| unique | 88 | 88 | 7 | 9 | 6 | 59 | 40 | 56 | 34 | 45 | 27 | 21 |
| top | ISO3 | Countries and areas | SSA | LAC | Upper middle income (UM) | 3% | 1% | 52% | 0% | 99% | Multiple Indicator Cluster Survey | 2018 |
| freq | 1 | 1 | 31 | 20 | 32 | 5 | 12 | 4 | 24 | 6 | 45 | 16 |

```
In [27]:    1  duplicates3 = df3.duplicated()
            2  print(f"Number of duplicate rows: {duplicates3.sum()}")
```

```
Number of duplicate rows: 0
```

```
In [28]:    1  df_transposed3 = df3.T
            2  duplicates3 = df_transposed3.duplicated()
            3  print(f"Number of duplicate columns: {duplicates3.sum()}")
            4  duplicate_column_names3 = df_transposed3.index[duplicates3].tolist()
            5  print("Duplicate column names:", duplicate_column_names3)
```

```
Number of duplicate columns: 0
Duplicate column names: []
```

```
In [29]:    1  df3.isnull().any()
```

```
Out[29]:  column A     False
          column B     False
          column C     False
          column  D    False
          column E     False
          column F     False
          column G      True
          column H      True
          column I      True
          column J      True
          column K     False
          column L     False
          dtype: bool
```

```
In [30]:    1  df3.isnull().sum()
```

```
Out[30]:  column A      0
          column B      0
          column C      0
          column  D     0
          column E      0
          column F      0
          column G     10
          column H      7
          column I     17
          column J     18
          column K      0
          column L      0
          dtype: int64
```

```
In [31]:   1  #removing the first 'column' row
           2  file_path_test3 = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSci
           3  dftest3 = pd.read_csv(file_path_test3, skiprows=1)
           4  output_file_path_test3 = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_
           5  dftest3.to_csv(output_file_path_test3, index=False)
           6  dftest3.head()
```

Out[31]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Time period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | DZA | Algeria | MENA | MENA | Upper middle income (UM) | 24% | 9% | 32% | 1% | 77% | Multiple Indicator Cluster Survey | 2018-19 |
| **1** | AGO | Angola | SSA | ESA | Lower middle income (LM) | 17% | 2% | 24% | 0% | 62% | Demographic and Health Survey | 2015-16 |
| **2** | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 40% | NaN | NaN | NaN | NaN | Multiple Indicator Cluster Survey | 2011-12 |
| **3** | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 81% | 71% | 88% | 47% | 99% | Demographic and Health Survey | 2015-16 |
| **4** | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 37% | 33% | 52% | 9% | 76% | Multiple Indicator Cluster Survey | 2019 |

```
In [32]:   1  #printing out every single row with missing values in the School Age Dataset
           2  missing_values_df_test3 = dftest3.loc[dftest3['Rural (Residence)'].isnull() | dftest3
```

```
In [33]:    1 missing_values_df_test3.head(100)
```

`Out[33]:`

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Time period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **2** | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 40% | NaN | NaN | NaN | NaN | Multiple Indicator Cluster Survey | 2011 1 |
| **7** | BOL | Bolivia (Plurinational State of) | LAC | LAC | Lower middle income (LM) | 12% | 4% | 17% | NaN | NaN | EDSA | 201 |
| **16** | CHN | China | EAP | EAP | Upper middle income (UM) | 57% | 50% | 91% | NaN | 90% | CHARLS | 201 |
| **24** | ECU | Ecuador | LAC | LAC | Upper middle income (UM) | 42% | 20% | 53% | NaN | NaN | ENSANUT | 201 |
| **25** | EGY | Egypt | MENA | MENA | Lower middle income (LM) | 17% | 9% | 29% | NaN | NaN | 2015 Household Income, Expenditure and Consump... | 201 |
| **37** | KEN | Kenya | SSA | ESA | Lower middle income (LM) | 32% | NaN | 24% | NaN | NaN | STEP Skills Measurement Household Survey 2013 ... | 201 |
| **46** | MEX | Mexico | LAC | LAC | Upper middle income (UM) | 41% | 11% | 52% | NaN | NaN | ENSANUT | 201 |
| **49** | MAR | Morocco | MENA | MENA | Lower middle income (LM) | 18% | 12% | 23% | NaN | NaN | Morocco Household and Youth Survey 2010 | 201 |
| **52** | NIC | Nicaragua | LAC | LAC | Lower middle income (LM) | 4% | NaN | NaN | NaN | NaN | Nicaragua National Demographic and Health Surv... | 2011 1 |
| **53** | NER | Niger | SSA | WCA | Low income (L) | 3% | NaN | NaN | NaN | NaN | National Survey on Household Living Conditions... | 2014 1 |
| **54** | NGA | Nigeria | SSA | WCA | Lower middle income (LM) | 3% | NaN | NaN | NaN | NaN | General Household Survey, Panel 2018-2019, Wave 4 | 2018 201 |
| **57** | PER | Peru | LAC | LAC | Upper middle income (UM) | 26% | 1% | 34% | NaN | NaN | ENDES | 201 |
| **59** | LCA | Saint Lucia | LAC | LAC | Upper middle income (UM) | 48% | 48% | 44% | 12% | NaN | Multiple Indicator Cluster Survey | 201 |
| **64** | SOM | Somalia | SSA | ESA | Low income (L) | 13% | NaN | NaN | NaN | NaN | Somalia High Frequency Survey | 2017 1 |

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Tim perio |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 65 | ZAF | South Africa | SSA | ESA | Upper middle income (UM) | 20% | NaN | NaN | NaN | NaN | South Africa Living Conditions Survey 2014-15 | 2014 1 |
| 80 | UKR | Ukraine | ECA | EECA | Lower middle income (LM) | 71% | NaN | 71% | NaN | NaN | STEP Skills Measurement Household Survey 2013 ... | 201 |
| 81 | GBR | United Kingdom | ECA | WE | High income (H) | 99% | NaN | NaN | NaN | NaN | UK Data Archive Information for the Study 8298... | 201 |
| 82 | URY | Uruguay | LAC | LAC | High income (H) | 63% | 47% | 65% | 35% | NaN | Multiple Indicator Cluster Survey | 2012 9 |
| 84 | VNM | Viet Nam | EAP | EAP | Lower middle income (LM) | 62% | NaN | 62% | NaN | NaN | STEP Skills Measurement Household Survey 2012 ... | 201 |

```
In [34]:  1  #I am going to be choosing Rural (Residence), and Urban (Residence) for analysis in
          2  #for this, like I did for task 1, I will be converting the % string values to float
          3
          4  def percentage_to_float(value):
          5      if isinstance(value, str) and value.endswith('%'):
          6          return float(value[:-1]) / 100
          7      else:
          8          return value
          9
         10  columns_to_convert = ['Total', 'Rural (Residence)', 'Urban (Residence)']
         11
         12  for column in columns_to_convert:
         13      dftest3[column] = dftest3[column].apply(percentage_to_float)
         14
         15  #removing unwanted last four columns
         16  #columns_to_remove = ['Poorest (Wealth quintile)', 'Richest (Wealth quintile)', 'Dat
         17  #dftest3 = dftest3.drop(columns=columns_to_remove)
         18
         19  dftest3.head()
```

Out[34]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Time period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | DZA | Algeria | MENA | MENA | Upper middle income (UM) | 0.24 | 0.09 | 0.32 | 1% | 77% | Multiple Indicator Cluster Survey | 2018-19 |
| 1 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.17 | 0.02 | 0.24 | 0% | 62% | Demographic and Health Survey | 2015-16 |
| 2 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.40 | NaN | NaN | NaN | NaN | Multiple Indicator Cluster Survey | 2011-12 |
| 3 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 | 0.71 | 0.88 | 47% | 99% | Demographic and Health Survey | 2015-16 |
| 4 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.37 | 0.33 | 0.52 | 9% | 76% | Multiple Indicator Cluster Survey | 2019 |

```
In [35]:  1  columns_to_remove = ['Poorest (Wealth quintile)', 'Richest (Wealth quintile)', 'Data
          2  dftest3 = dftest3.drop(columns=columns_to_remove)
          3  dftest3.head()
```

Out[35]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) |
|---|---|---|---|---|---|---|---|---|
| 0 | DZA | Algeria | MENA | MENA | Upper middle income (UM) | 0.24 | 0.09 | 0.32 |
| 1 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.17 | 0.02 | 0.24 |
| 2 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.40 | NaN | NaN |
| 3 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 | 0.71 | 0.88 |
| 4 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.37 | 0.33 | 0.52 |

```
In [36]:    1  file_path = "/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSciWithP
            2  dftest3.to_csv(file_path, index=False)
```

```
In [37]:    1
            2  file_path = "/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSciWithP
            3  dftest3cleaned1 = pd.read_csv(file_path)
            4
            5  dftest3cleaned1.head()
```

Out[37]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) |
|---|---|---|---|---|---|---|---|---|
| 0 | DZA | Algeria | MENA | MENA | Upper middle income (UM) | 0.24 | 0.09 | 0.32 |
| 1 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.17 | 0.02 | 0.24 |
| 2 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.40 | NaN | NaN |
| 3 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 | 0.71 | 0.88 |
| 4 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.37 | 0.33 | 0.52 |

```
In [38]:   1
           2   # Calculating the mean values for 'Rural (Residence)' and 'Urban (Residence)' per Reg
           3   mean_values = dftest3cleaned1.groupby('Region').agg({
           4       'Rural (Residence)': np.nanmean,
           5       'Urban (Residence)': np.nanmean
           6   }).reset_index()
           7
           8
           9   # Function to substitute NaN values with the mean value based on the Region's mean
          10   def replace_nan_with_mean(row, mean_values):
          11       region = row['Region']
          12       mean_rural = mean_values.loc[mean_values['Region'] == region, 'Rural (Residence)
          13       mean_urban = mean_values.loc[mean_values['Region'] == region, 'Urban (Residence)
          14
          15       row['Rural (Residence)'] = round(row['Rural (Residence)'] if pd.notna(row['Rural
          16       row['Urban (Residence)'] = round(row['Urban (Residence)'] if pd.notna(row['Urban
          17
          18       return row
          19
          20
          21   # Substituting NaN values into dftest3cleaned1
          22   dftest3cleaned1 = dftest3cleaned1.apply(lambda row: replace_nan_with_mean(row, mean_v
          23
          24   dftest3cleaned1.head()
          25
```

Out[38]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) |
|---|---|---|---|---|---|---|---|---|
| 0 | DZA | Algeria | MENA | MENA | Upper middle income (UM) | 0.24 | 0.09 | 0.32 |
| 1 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.17 | 0.02 | 0.24 |
| 2 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.40 | 0.27 | 0.49 |
| 3 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 | 0.71 | 0.88 |
| 4 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.37 | 0.33 | 0.52 |

```
In [39]:   1  dftest3cleaned1.head(50)
```

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) |
|---|---|---|---|---|---|---|---|---|
| 0 | DZA | Algeria | MENA | MENA | Upper middle income (UM) | 0.24 | 0.09 | 0.32 |
| 1 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.17 | 0.02 | 0.24 |
| 2 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.40 | 0.27 | 0.49 |
| 3 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 | 0.71 | 0.88 |
| 4 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.37 | 0.33 | 0.52 |
| 5 | BRB | Barbados | LAC | LAC | High income (H) | 0.66 | 0.61 | 0.69 |
| 6 | BEN | Benin | SSA | WCA | Low income (L) | 0.04 | 0.01 | 0.07 |
| 7 | BOL | Bolivia (Plurinational State of) | LAC | LAC | Lower middle income (LM) | 0.12 | 0.04 | 0.17 |
| 8 | BIH | Bosnia and Herzegovina | ECA | EECA | Upper middle income (UM) | 0.59 | 0.51 | 0.76 |
| 9 | BRA | Brazil | LAC | LAC | Upper middle income (UM) | 0.83 | 0.51 | 0.89 |
| 10 | BGR | Bulgaria | ECA | EECA | Upper middle income (UM) | 0.76 | 0.66 | 0.81 |
| 11 | BFA | Burkina Faso | SSA | WCA | Low income (L) | 0.01 | 0.01 | 0.04 |
| 12 | CMR | Cameroon | SSA | WCA | Lower middle income (LM) | 0.05 | 0.00 | 0.10 |
| 13 | CAF | Central African Republic | SSA | WCA | Low income (L) | 0.04 | 0.01 | 0.09 |
| 14 | TCD | Chad | SSA | WCA | Low income (L) | 0.02 | 0.01 | 0.08 |
| 15 | CHL | Chile | LAC | LAC | High income (H) | 0.86 | 0.70 | 0.89 |
| 16 | CHN | China | EAP | EAP | Upper middle income (UM) | 0.57 | 0.50 | 0.91 |
| 17 | COL | Colombia | LAC | LAC | Upper middle income (UM) | 0.36 | 0.05 | 0.48 |
| 18 | CRI | Costa Rica | LAC | LAC | Upper middle income (UM) | 0.72 | 0.61 | 0.78 |
| 19 | CIV | C™te d'Ivoire | SSA | WCA | Lower middle income (LM) | 0.03 | 0.01 | 0.05 |
| 20 | CUB | Cuba | LAC | LAC | Upper middle income (UM) | 0.04 | 0.01 | 0.06 |
| 21 | COD | Democratic Republic of the Congo | SSA | WCA | Low income (L) | 0.01 | 0.00 | 0.02 |
| 22 | DJI | Djibouti | SSA | ESA | Lower middle income (LM) | 0.06 | 0.01 | 0.11 |
| 23 | DOM | Dominican Republic | LAC | LAC | Upper middle income (UM) | 0.24 | 0.10 | 0.29 |
| 24 | ECU | Ecuador | LAC | LAC | Upper middle income (UM) | 0.42 | 0.20 | 0.53 |
| 25 | EGY | Egypt | MENA | MENA | Lower middle income (LM) | 0.17 | 0.09 | 0.29 |
| 26 | GMB | Gambia | SSA | WCA | Low income (L) | 0.65 | 0.47 | 0.75 |
| 27 | GEO | Georgia | ECA | EECA | Upper middle income (UM) | 0.85 | 0.72 | 0.93 |
| 28 | GHA | Ghana | SSA | WCA | Lower middle income (LM) | 0.17 | 0.10 | 0.26 |

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) |
|---|---|---|---|---|---|---|---|---|
| 29 | GTM | Guatemala | LAC | LAC | Upper middle income (UM) | 0.09 | 0.03 | 0.18 |
| 30 | GNB | Guinea-Bissau | SSA | WCA | Low income (L) | 0.02 | 0.02 | 0.04 |
| 31 | HTI | Haiti | LAC | LAC | Low income (L) | 0.21 | 0.11 | 0.36 |
| 32 | IND | India | SA | SA | Lower middle income (LM) | 0.09 | 0.05 | 0.17 |
| 33 | IDN | Indonesia | EAP | EAP | Lower middle income (LM) | 0.19 | 0.10 | 0.26 |
| 34 | IRQ | Iraq | MENA | MENA | Upper middle income (UM) | 0.49 | 0.35 | 0.56 |
| 35 | JPN | Japan | EAP | EAP | High income (H) | 0.78 | 0.83 | 0.77 |
| 36 | JOR | Jordan | MENA | MENA | Upper middle income (UM) | 0.38 | 0.35 | 0.38 |
| 37 | KEN | Kenya | SSA | ESA | Lower middle income (LM) | 0.32 | 0.07 | 0.24 |
| 38 | KIR | Kiribati | EAP | EAP | Lower middle income (LM) | 0.51 | 0.33 | 0.69 |
| 39 | KGZ | Kyrgyzstan | ECA | EECA | Lower middle income (LM) | 0.74 | 0.69 | 0.83 |
| 40 | LAO | Lao People's Democratic Republic | EAP | EAP | Lower middle income (LM) | 0.02 | 0.01 | 0.04 |
| 41 | LSO | Lesotho | SSA | ESA | Lower middle income (LM) | 0.32 | 0.22 | 0.52 |
| 42 | MDG | Madagascar | SSA | ESA | Low income (L) | 0.11 | 0.06 | 0.27 |
| 43 | MDV | Maldives | SA | SA | Upper middle income (UM) | 0.70 | 0.69 | 0.72 |
| 44 | MLI | Mali | SSA | WCA | Low income (L) | 0.05 | 0.03 | 0.14 |
| 45 | MRT | Mauritania | SSA | WCA | Lower middle income (LM) | 0.03 | 0.01 | 0.06 |
| 46 | MEX | Mexico | LAC | LAC | Upper middle income (UM) | 0.41 | 0.11 | 0.52 |
| 47 | MNG | Mongolia | EAP | EAP | Lower middle income (LM) | 0.37 | 0.13 | 0.49 |
| 48 | MNE | Montenegro | ECA | EECA | Upper middle income (UM) | 0.82 | 0.74 | 0.86 |
| 49 | MAR | Morocco | MENA | MENA | Lower middle income (LM) | 0.18 | 0.12 | 0.23 |

```
In [40]:  1  file_path3 = "/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSciWith
          2  dftest3cleaned1.to_csv(file_path3, index=False)
```

```
1
2  file_path = "/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSciWithPy
3  dftest3cleaned2 = pd.read_csv(file_path)
4
5  dftest3cleaned2.head()
```

Out[41]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) |
|---|---|---|---|---|---|---|---|---|
| 0 | DZA | Algeria | MENA | MENA | Upper middle income (UM) | 0.24 | 0.09 | 0.32 |
| 1 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.17 | 0.02 | 0.24 |
| 2 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.40 | 0.27 | 0.49 |
| 3 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 | 0.71 | 0.88 |
| 4 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.37 | 0.33 | 0.52 |

## Task 2.2: Top 10 Countries and areas with the Highest percent of School-aged children with internet access

This is the start of Task 2.2's objectives

```
In [42]:  1
          2
          3  dftest3cleaned2 = dftest3cleaned2.sort_values(by='Total', ascending=False)
          4
          5  #selecting top 10 rows based on their 'Total' value
          6  top_10 = dftest3cleaned2.head(10)
          7
          8  # Bar graph for displaying the top 10 countries with highest total percent of school
          9  plt.figure(figsize=(12, 6))
         10  plt.bar(top_10['Countries and areas'], top_10['Total'])
         11  plt.xlabel('Countries and areas')
         12  plt.ylabel('Total')
         13  plt.title('Top 10 Countries and Areas with the Highest percent of School-aged childre
         14  plt.xticks(rotation=45)
         15
         16  plt.tight_layout()
         17  plt.show()
         18
```



Top 10 Countries and Areas with the Highest percent of School-aged children with internet access

```
In [43]:  1  #horizintal version of the above graph
          2  plt.figure(figsize=(12, 6))
          3  plt.barh(top_10['Countries and areas'], top_10['Total'])
          4  plt.xlabel('Total')
          5  plt.ylabel('Countries and areas')
          6  plt.title('Top 10 Countries and Areas with the Highest Total(%) of School-aged childr
          7  plt.tight_layout()
          8  plt.show()
          9
```

Top 10 Countries and Areas with the Highest Total(%) of School-aged children with internet access

```
In [44]:   1  #here, I am visualizing the top 10 countries which I pulled into a dataframe
           2  #and depictted in a bar graph above, but am delineating by the Top 10's Income Group:
           3  #to see the relationship between Income group and the rural versus urban residence
           4  #in internet connectivitty
           5
           6
           7  dftest3cleaned2 = dftest3cleaned2.sort_values(by='Total', ascending=False)
           8
           9  top_10 = dftest3cleaned2.head(10)
          10
          11  fig, axes = plt.subplots(2, 2, figsize=(12, 12))
          12  income_groups = ['Upper middle income (UM)', 'Lower middle income (LM)', 'High income
          13
          14  for i, income_group in enumerate(income_groups):
          15      row = i // 2
          16      col = i % 2
          17
          18      top_10_income_group = top_10[top_10['Income Group'] == income_group]
          19
          20      if not top_10_income_group.empty:
          21          ax = sns.boxplot(data=top_10_income_group[['Rural (Residence)', 'Urban (Resid
          22          ax.set_xticklabels(['Rural (Residence)', 'Urban (Residence)'])
          23          ax.set_ylabel('Total percent')
          24          ax.set_ylim(0, 1)
          25          ax.yaxis.grid(False)
          26          axes[row, col].set_title(f'Top 10 countries and areas with the highest total
          27      else:
          28          axes[row, col].set_axis_off()
          29
          30  plt.subplots_adjust(hspace=0.3)
          31  plt.show()
          32
```

Top 10 countries and areas with the highest total percentage of school-age children in terms of their Income Group and Residence
Income Group: Upper middle income (UM)



Top 10 countries and areas with the highest total percentage of school-age children in terms of their Income Group and Residence
Income Group: High income (H)

```
In [45]:   1  #here, i am separating the bar graphs to be able to visualize the two
           2  #relevant income groups more effectiveely I will create a groupeed Bar plot
           3  #these are visually quite explainable to a viewer and I think that makes it a favoral
           4  #visualization method
           5
           6
           7  #melting the data to fit into a grouped bar plot
           8  melted_data = pd.melt(top_10, id_vars=['Income Group'], value_vars=['Rural (Residence
           9
          10  plt.figure(figsize=(12, 6))
          11  ax = sns.barplot(data=melted_data, x='Income Group', y='Total percent', hue='Residence
          12  plt.title('Top 10 countries and areas with the highest total percentage of school-age
          13  plt.ylabel('Total percent')
          14
          15
          16  ax.set_ylim(0, 1)
          17
          18  plt.show()
          19
```

Top 10 countries and areas with the highest total percentage of school-age children in terms of their Income Group and Residence

```
In [46]:   1
           2  # Here I am Melting the data to have the columns, 'Residence' and 'Total percent'
           3  #however, i am going to be visualizing how many individual values we have per catego
           4  #this is important as it gives a better understanding of what kind of data visualiza
           5  melted_data = pd.melt(top_10, id_vars=['Income Group'], value_vars=['Rural (Residence
           6
           7  # Faceted Histogram will help me in grasping how many individual values are present
           8  #for me to be able to visualize and make inferences with
           9  g = sns.FacetGrid(melted_data, col='Income Group', hue='Residence', palette='Set2',
          10  g.map(sns.histplot, 'Total percent', bins=10, alpha=0.75)
          11  g.add_legend(title='Residence')
          12  g.set_axis_labels('Total percent', 'Frequency')
          13  g.fig.subplots_adjust(wspace=0.3, hspace=0.3)
          14  plt.show()
          15
```

```
In [47]:   #Here I want to visualize the relationship between the urban and rural residence
           #with a Scatter plot
        3
           plt.figure(figsize=(12, 6))
        5
           #First I am going to define markers for each income group
           markers = {'Upper middle income (UM)': 'o', 'Lower middle income (LM)': 's', 'High inco
        8
        9
           sns.scatterplot(x='Rural (Residence)', y='Urban (Residence)', hue='Income Group', data=
       11
           # This regression line will help in visualizing the relationship between urban and rura
           sns.regplot(x='Rural (Residence)', y='Urban (Residence)', data=top_10, scatter=False, l
       14
           plt.title('Scatter Plot with Regression Line of Rural and Urban Residence')
           plt.xlabel('Rural (Residence)')
           plt.ylabel('Urban (Residence)')
       18
           # Here, i am setting the y-axis range from 0 to 1 so that we see the clear spread acros
           #percentages, in the next cell, I will zoom into the relevant section of this scatter p
           #examine the relationship in more detail
           plt.ylim(0, 1)
           plt.xlim(0,1)
       24
           plt.legend(title='Income Group', loc='upper left')
       26
           plt.show()
       28
```



Scatter Plot with Regression Line of Rural and Urban Residence

```
In [48]:   1  #zoomed in scatter plot
           2  #given that the previous scatter plot gave a birds eye view of the data, It prompts
           3  #a zoomed in view of the data as well
           4
           5  plt.figure(figsize=(12, 6))
           6
           7  #First I am going to define markers for each income group
           8  markers = {'Upper middle income (UM)': 'o', 'Lower middle income (LM)': 's', 'High in
           9
          10  sns.scatterplot(x='Rural (Residence)', y='Urban (Residence)', hue='Income Group', dat
          11
          12  # This regression line will help in visualizing the relationship between urban and ru
          13  #given that this graph will be zoomed in, it will be easier to make inferences
          14  sns.regplot(x='Rural (Residence)', y='Urban (Residence)', data=top_10, scatter=False
          15
          16  plt.title('Relationship Between percent of Urban and Rural School Age Children with
          17  plt.xlabel('Rural (Residence)')
          18  plt.ylabel('Urban (Residence)')
          19
          20  plt.legend(title='Income Group', loc='upper left')
          21
          22  plt.show()
          23
```



Relationship Between percent of Urban and Rural School Age Children with Internet Connection

## Task 2.3: Preparing Secondary.csv data for comparative analysis with cleaned Primary.csv data

In this part, I am cleaning Secondary.csv's data and handling the anomalies so that I can compare the results to the Primary.csv data to see if the results mirror their respective regions

```
In [49]:  1  #cleaning Secondary.csv
          2  file_path_test2 = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataSc
          3  dftest2 = pd.read_csv(file_path_test2, skiprows=1)
          4  output_file_path_test2 = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_
          5  dftest2.to_csv(output_file_path_test2, index=False)
          6  dftest2.head()
```

Out[49]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Time period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | AGOA | Angola | SSA | ESA | Lower middle income (LM) | 24% | 2% | 33% | 0% | 69% | Demographic and Health Survey | 2015-16 |
| 1 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 45% | NaN | NaN | NaN | NaN | Multiple Indicator Cluster Survey | 2011-12 |
| 2 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 85% | 78% | 91% | 54% | 100% | Demographic and Health Survey | 2015-16 |
| 3 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 42% | 38% | 57% | 13% | 79% | Multiple Indicator Cluster Survey | 2019 |
| 4 | BRB | Barbados | LAC | LAC | High income (H) | 76% | 76% | 76% | 4% | 100% | Multiple Indicator Cluster Survey | 2012 |

```
In [50]:  1  #converting the data in the 'Total' column to numeric to enable analysis
          2
          3  def percentage_to_float(value):
          4      if isinstance(value, str) and value.endswith('%'):
          5          return float(value[:-1]) / 100
          6      else:
          7          return value
          8
          9  columns_to_convert = ['Total']
         10
         11  for column in columns_to_convert:
         12      dftest2[column] = dftest2[column].apply(percentage_to_float)
         13
         14  dftest2.head()
```

Out[50]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total | Rural (Residence) | Urban (Residence) | Poorest (Wealth quintile) | Richest (Wealth quintile) | Data source | Time period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | AGOA | Angola | SSA | ESA | Lower middle income (LM) | 0.24 | 2% | 33% | 0% | 69% | Demographic and Health Survey | 2015-16 |
| **1** | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.45 | NaN | NaN | NaN | NaN | Multiple Indicator Cluster Survey | 2011-12 |
| **2** | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.85 | 78% | 91% | 54% | 100% | Demographic and Health Survey | 2015-16 |
| **3** | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.42 | 38% | 57% | 13% | 79% | Multiple Indicator Cluster Survey | 2019 |
| **4** | BRB | Barbados | LAC | LAC | High income (H) | 0.76 | 76% | 76% | 4% | 100% | Multiple Indicator Cluster Survey | 2012 |

```
In [51]:  1  columns_to_remove2 = ['Rural (Residence)','Urban (Residence)','Poorest (Wealth quintile)
          2  dftest2 = dftest2.drop(columns=columns_to_remove2)
          3  dftest2.head()
```

Out[51]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|---|---|---|---|---|---|---|
| **0** | AGOA | Angola | SSA | ESA | Lower middle income (LM) | 0.24 |
| **1** | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.45 |
| **2** | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.85 |
| **3** | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.42 |
| **4** | BRB | Barbados | LAC | LAC | High income (H) | 0.76 |

```
In [52]:  1  output_file_path2 = '/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignment1_DataS
          2  dftest2.to_csv(output_file_path2, index=False)
```

```
In [53]:   1  dftest.head()
```

Out[53]:

|   | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|---|------|---------------------|--------|------------|--------------|-------|
| 0 | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.15 |
| 1 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.39 |
| 2 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.81 |
| 3 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.34 |
| 4 | BRB | Barbados | LAC | LAC | High income (H) | 0.63 |

```
In [54]:   1  dftest2.head()
```

Out[54]:

|   | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|---|------|---------------------|--------|------------|--------------|-------|
| 0 | AGOA | Angola | SSA | ESA | Lower middle income (LM) | 0.24 |
| 1 | ARG | Argentina | LAC | LAC | Upper middle income (UM) | 0.45 |
| 2 | ARM | Armenia | ECA | EECA | Upper middle income (UM) | 0.85 |
| 3 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.42 |
| 4 | BRB | Barbados | LAC | LAC | High income (H) | 0.76 |

```
In [55]:   1  dftesttask3 = dftest[dftest['Income Group'] == 'Lower middle income (LM)']
           2  dftesttask3.head()
```

Out[55]:

|    | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|----|------|---------------------|--------|------------|--------------|-------|
| 0  | AGO | Angola | SSA | ESA | Lower middle income (LM) | 0.15 |
| 3  | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.34 |
| 6  | BOL | Bolivia (Plurinational State of) | LAC | LAC | Lower middle income (LM) | 0.11 |
| 11 | CMR | Cameroon | SSA | WCA | Lower middle income (LM) | 0.04 |
| 18 | CIV | C™te d'Ivoire | SSA | WCA | Lower middle income (LM) | 0.02 |

```
In [56]:   1  dftest2task3 = dftest2[dftest2['Income Group'] == 'Lower middle income (LM)']
           2
           3  dftest2task3.head()
```

Out[56]:

|    | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|----|------|---------------------|--------|------------|--------------|-------|
| 0  | AGOA | Angola | SSA | ESA | Lower middle income (LM) | 0.24 |
| 3  | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.42 |
| 10 | CMR | Cameroon | SSA | WCA | Lower middle income (LM) | 0.07 |
| 17 | CIV | C™te d'Ivoire | SSA | WCA | Lower middle income (LM) | 0.03 |
| 20 | DJI | Djibouti | SSA | ESA | Lower middle income (LM) | 0.09 |

```
In [57]:   1  file_path_primarylm_task3 = "/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignmen
           2
           3  dftesttask3.to_csv(file_path_primarylm_task3, index=False)
           4
```

```
In [58]:    1  file_path_secondarylm_task3 = "/Users/amayiyer/Desktop/DatSci_Python/s3970066/Assignm
            2
            3  dftest2task3.to_csv(file_path_secondarylm_task3, index=False)
```

```
In [59]:    1  plt.figure(figsize=(12, 6))
            2  sns.boxplot(x='Region', y='Total', data=dftesttask3)
            3  plt.title('Box Plot of Primary Children with Internet Access for Lower Middle income
            4  plt.xlabel('Region')
            5  plt.ylabel('Total')
            6  plt.show()
```



Box Plot of Primary Children with Internet Access for Lower Middle income countries by Region

```
1
2  plt.figure(figsize=(12, 6))
3  sns.violinplot(x='Region', y='Total', data=dftesttask3)
4  plt.title('Violin Plot of Primary Children with Internet Access for Lower Middle Inc
5  plt.xlabel('Region')
6  plt.ylabel('Total')
7  plt.show()
8
```



Violin Plot of Primary Children with Internet Access for Lower Middle Income Countries by Region

```
In [61]:  1
          2  plt.figure(figsize=(12, 6))
          3  sns.swarmplot(x='Region', y='Total', data=dftesttask3)
          4  plt.title('Swarm Plot of Primary Children with Internet Access for Lower Middle Incor
          5  plt.xlabel('Region')
          6  plt.ylabel('Total')
          7  plt.show()
          8
```



Swarm Plot of Primary Children with Internet Access for Lower Middle Income Countries by Region

```
In [62]:  1
          2  plt.figure(figsize=(12, 6))
          3  sns.barplot(x='Region', y='Total', data=dftesttask3, ci='sd')
          4  plt.title('Bar Plot of Primary Children with Internet Access for Lower Middle Income
          5  plt.xlabel('Region')
          6  plt.ylabel('Total')
          7  plt.show()
          8
```



Bar Plot of Primary Children with Internet Access for Lower Middle Income Countries by Region

```
1
2  plt.figure(figsize=(12, 6))
3  sns.barplot(x='Region', y='Total', data=dftesttask3, ci='sd')
4  plt.title('Bar Plot of Primary Children with Internet Access for Lower Middle Income
5  plt.xlabel('Region')
6  plt.ylabel('Total')
7  plt.show()
8
```



Bar Plot of Primary Children with Internet Access for Lower Middle Income Countries by Region

```
1
2  plt.figure(figsize=(12, 6))
3  sns.scatterplot(x=dftesttask3.index, y='Total', hue='Region', data=dftesttask3)
4  plt.title('Scatter Plot of Primary Children with Internet Access for Lower Middle Inc
5  plt.xlabel('Index')
6  plt.ylabel('Total')
7  plt.show()
8
```

Scatter Plot of Primary Children with Internet Access for Lower Middle Income Countries by Region

```
1
2  plt.figure(figsize=(12, 6))
3  sns.stripplot(x='Region', y='Total', data=dftesttask3, jitter=True)
4  plt.title('Strip Plot of Primary Children with Internet Access for Lower Middle Incom
5  plt.xlabel('Region')
6  plt.ylabel('Total')
7  plt.show()
8
9
```
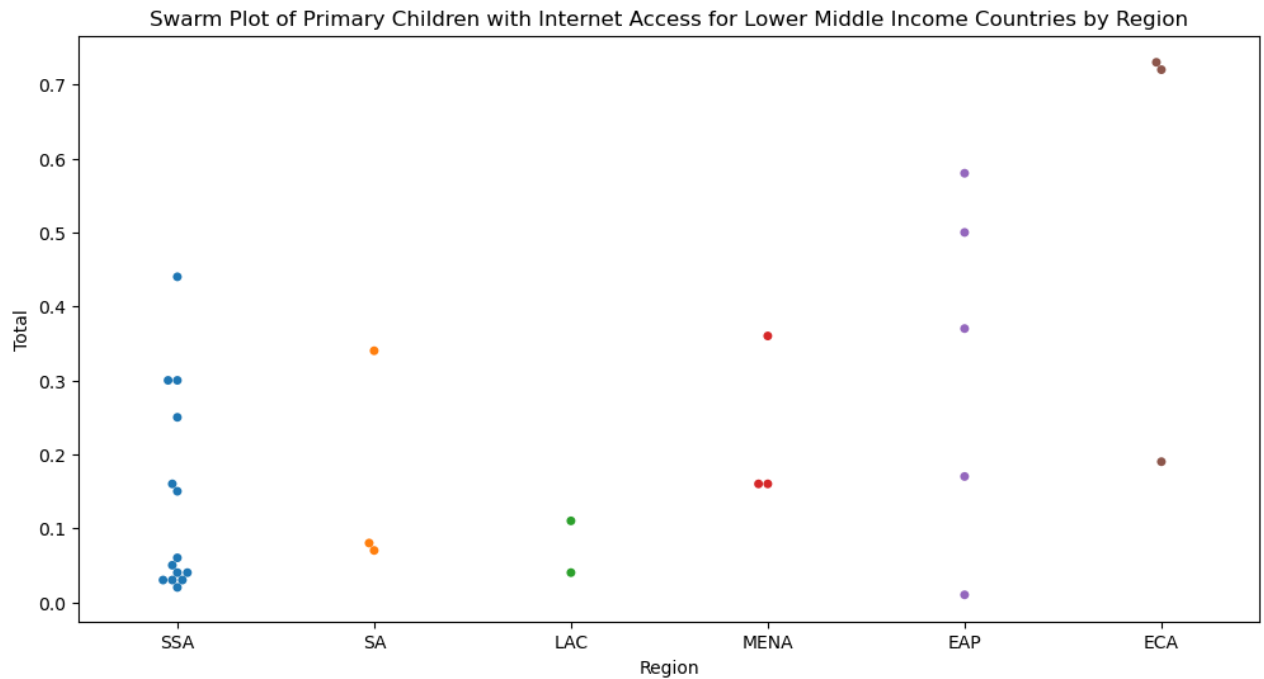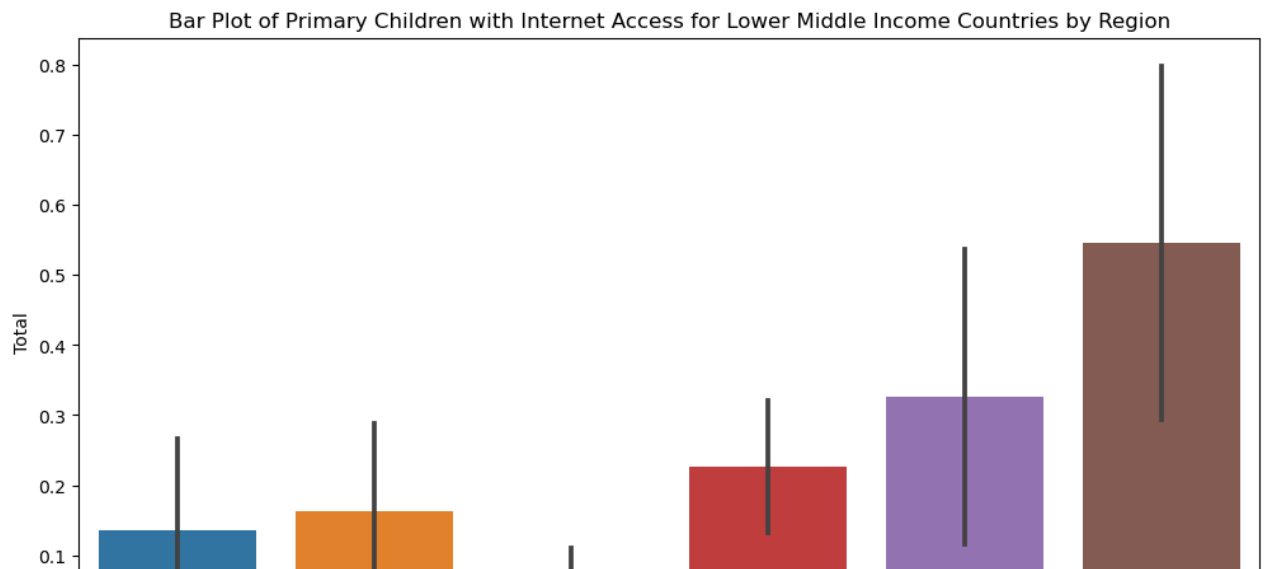


Strip Plot of Primary Children with Internet Access for Lower Middle Income Countries by Region

```
In [66]:   1  #I am visualizing the following data in a Faceted histogram to get an idea of how mai
           2  #regions in the datasets have a low % of internet connectivity
           3
           4  import seaborn as sns
           5  import matplotlib.pyplot as plt
           6
           7  g = sns.FacetGrid(dftesttask3, hue='Region', height=6, aspect=2)
           8  g.map(sns.histplot, 'Total', kde=True, bins=20)
           9  g.add_legend()
          10  plt.title('Histogram and KDE of Primary Children with Internet Access for Lower Middl
          11  plt.xlabel('Total')
          12  plt.show()
          13
```



Histogram and KDE of Primary Children with Internet Access for Lower Middle Income Countries by Region

```
1  plt.figure(figsize=(12, 6))
2  sns.boxplot(x='Region', y='Total', data=dftest2task3)
3  plt.title('Box Plot of Secondary Children with Internet Access for Lower Middle incor
4  plt.xlabel('Region')
5  plt.ylabel('Total')
6  plt.show()
```

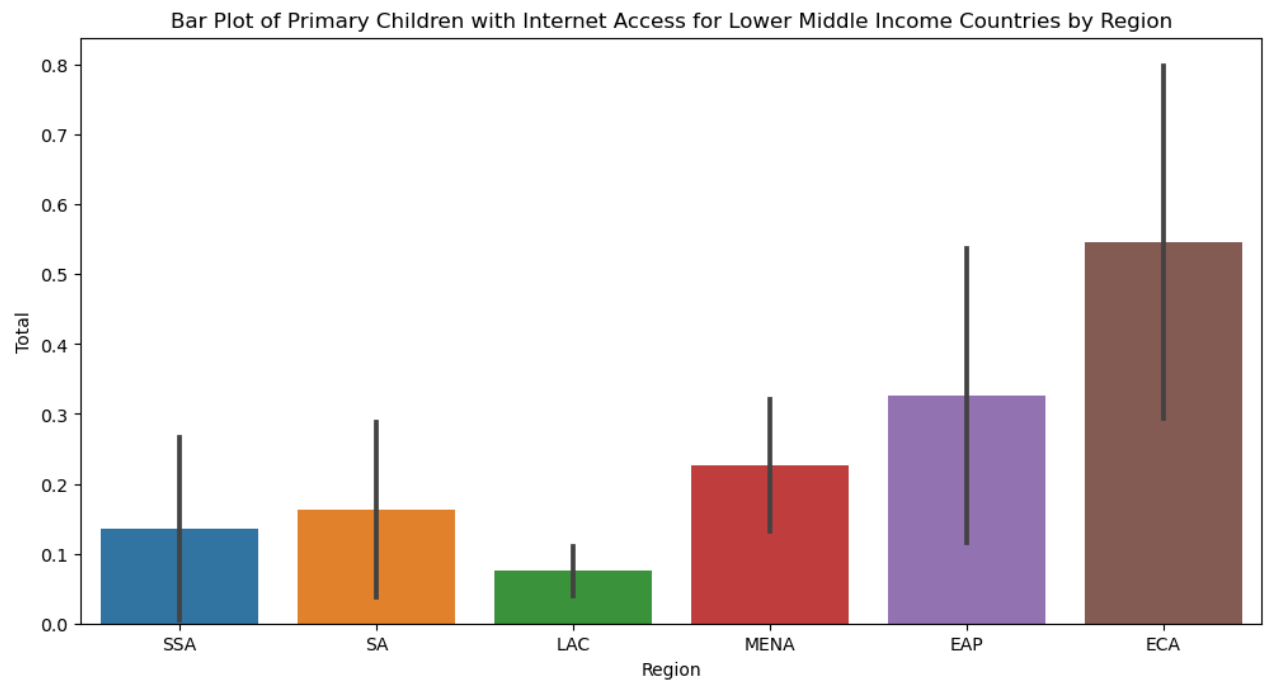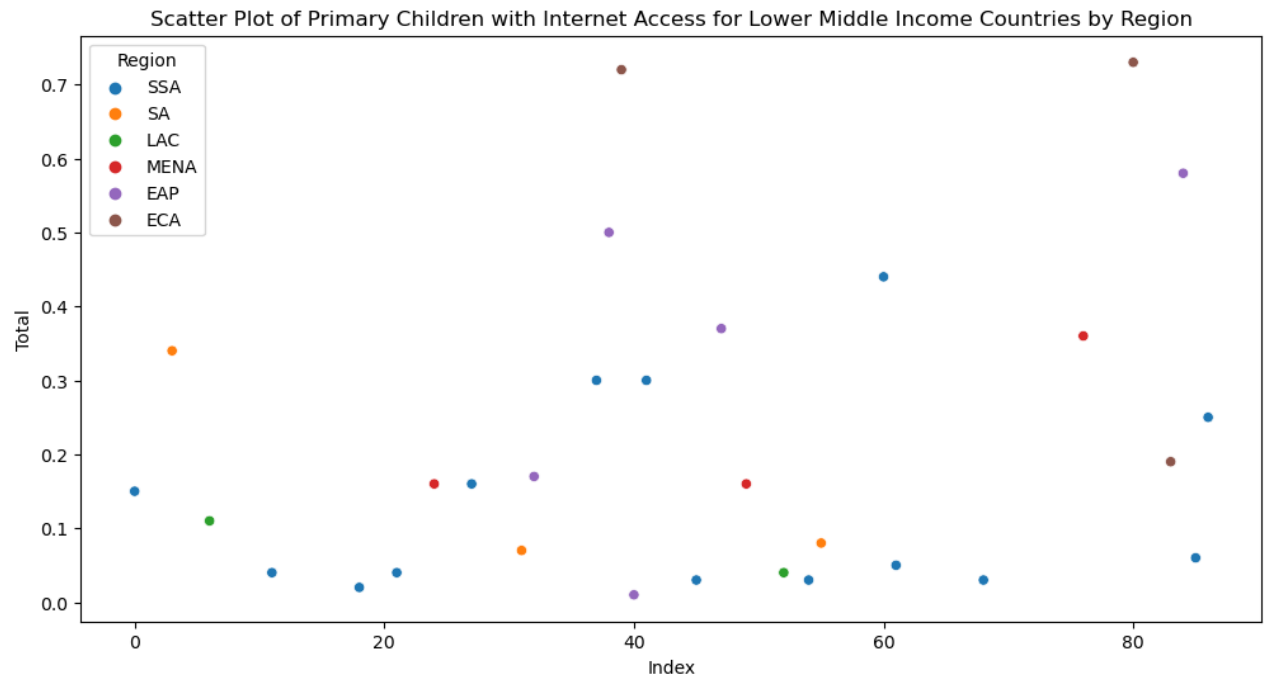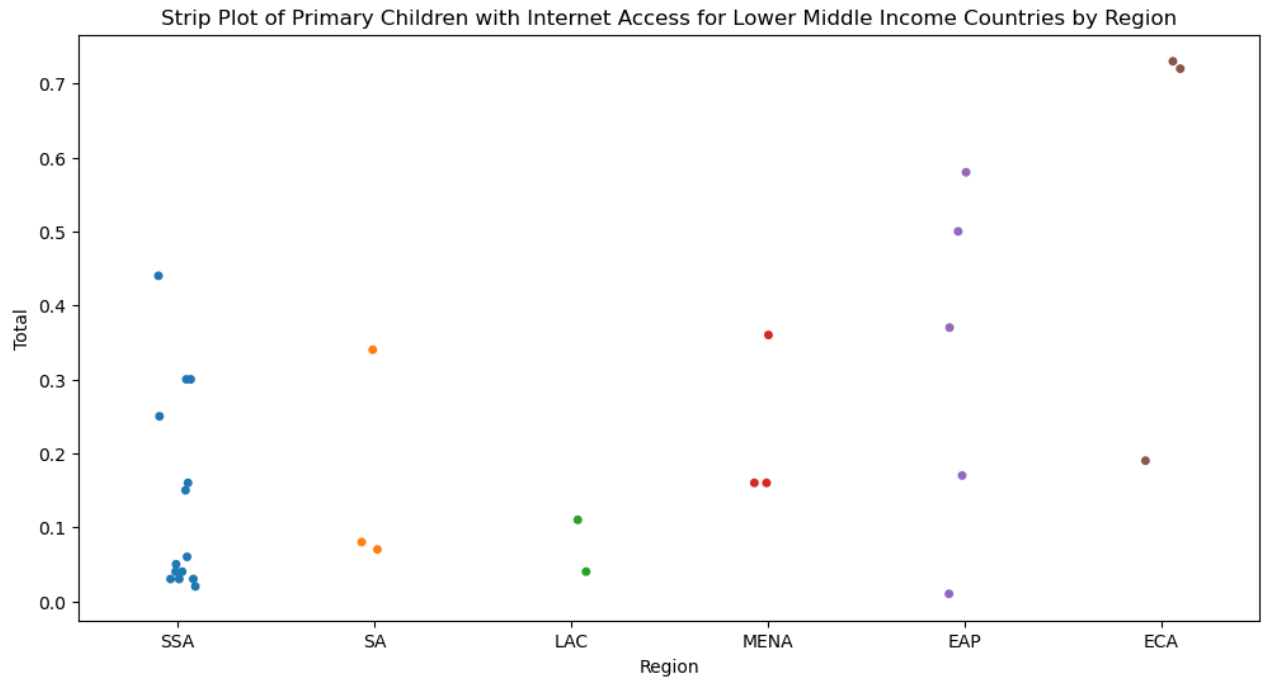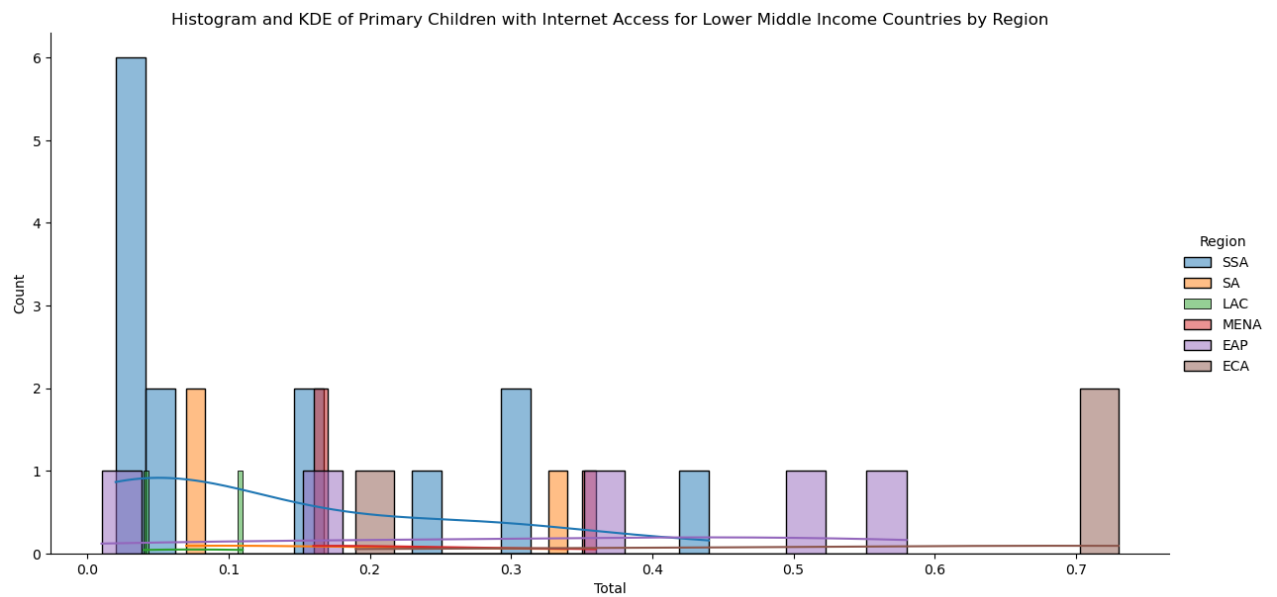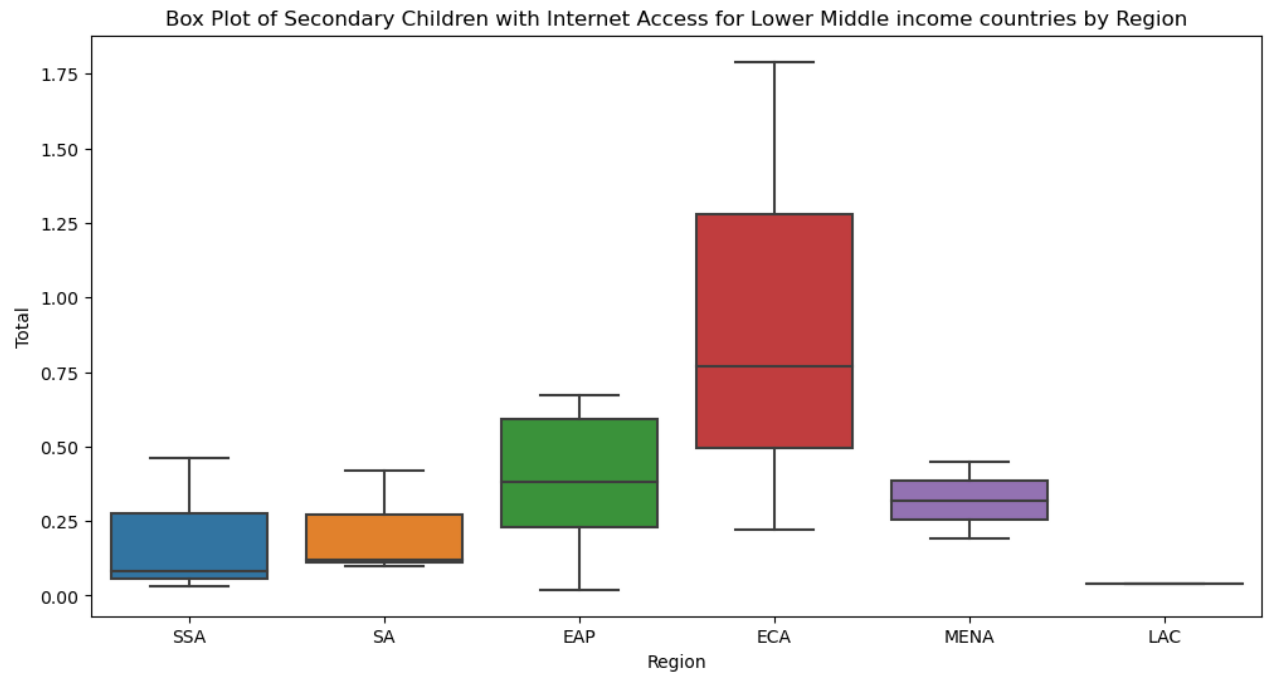Box Plot of Secondary Children with Internet Access for Lower Middle income countries by Region

```
In [68]:    1 dftest2task3.head(50)
```

Out[68]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|---|---|---|---|---|---|---|
| 0 | AGOA | Angola | SSA | ESA | Lower middle income (LM) | 0.24 |
| 3 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.42 |
| 10 | CMR | Cameroon | SSA | WCA | Lower middle income (LM) | 0.07 |
| 17 | CIV | C™te d'Ivoire | SSA | WCA | Lower middle income (LM) | 0.03 |
| 20 | DJI | Djibouti | SSA | ESA | Lower middle income (LM) | 0.09 |
| 24 | GHA | Ghana | SSA | WCA | Lower middle income (LM) | 0.20 |
| 28 | IND | India | SA | SA | Lower middle income (LM) | 0.12 |
| 29 | IDN | Indonesia | EAP | EAP | Lower middle income (LM) | 0.23 |
| 34 | KEN | Kenya | SSA | ESA | Lower middle income (LM) | 0.39 |
| 35 | KIR | Kiribati | EAP | EAP | Lower middle income (LM) | 0.59 |
| 36 | KGZ | Kyrgyzstan | ECA | EECA | Lower middle income (LM) | 0.77 |
| 37 | LAO | Lao People's Democratic Republic | EAP | EAP | Lower middle income (LM) | 0.02 |
| 38 | LSO | Lesotho | SSA | ESA | Lower middle income (LM) | 0.38 |
| 42 | MRT | Mauritania | SSA | WCA | Lower middle income (LM) | 0.04 |
| 43 | MNG | Mongolia | EAP | EAP | Lower middle income (LM) | 0.38 |
| 45 | MAR | Morocco | MENA | MENA | Lower middle income (LM) | 0.19 |
| 48 | NIC | Nicaragua | LAC | LAC | Lower middle income (LM) | 0.04 |
| 50 | NGA | Nigeria | SSA | WCA | Lower middle income (LM) | 0.04 |
| 51 | PAK | Pakistan | SA | SA | Lower middle income (LM) | 0.10 |
| 55 | STP | Sao Tome and Principe | SSA | WCA | Lower middle income (LM) | 0.46 |
| 56 | SEN | Senegal | SSA | WCA | Lower middle income (LM) | 0.07 |
| 63 | SDN | Sudan | SSA | ESA | Lower middle income (LM) | 0.05 |
| 71 | TUN | Tunisia | MENA | MENA | Lower middle income (LM) | 0.45 |
| 75 | UKR | Ukraine | ECA | EECA | Lower middle income (LM) | 1.79 |
| 78 | UZB | Uzbekistan | ECA | EECA | Lower middle income (LM) | 0.22 |
| 79 | VNM | Viet Nam | EAP | EAP | Lower middle income (LM) | 0.67 |
| 80 | ZMB | Zambia | SSA | ESA | Lower middle income (LM) | 0.07 |
| 81 | ZWE | Zimbabwe | SSA | ESA | Lower middle income (LM) | 0.29 |

```
In [69]:   1  #because of Ukraine's total value being over 1, I am handling the anomaly by
           2  #making it 0.79, along with this, I will also handle other similarly anomalous values
           3  #although Ukraine is the only main anomaly here
           4
           5  def adjust_total(value):
           6      if value < 0:
           7          return -value
           8      elif value > 1:
           9          return value - int(value)
          10      else:
          11          return value
          12
          13
          14  dftest2task3['Total'] = dftest2task3['Total'].apply(adjust_total)
          15
          16
          17  dftest2task3.head(50)
```

```
/var/folders/8y/9bzf6lrd1kggm5429dypyz2r0000gn/T/ipykernel_40815/3478238393.py:14: Set
tingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/use
r_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas
-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
  dftest2task3['Total'] = dftest2task3['Total'].apply(adjust_total)
```

Out[69]:

| | ISO3 | Countries and areas | Region | Sub-region | Income Group | Total |
|---|---|---|---|---|---|---|
| 0 | AGOA | Angola | SSA | ESA | Lower middle income (LM) | 0.24 |
| 3 | BGD | Bangladesh | SA | SA | Lower middle income (LM) | 0.42 |
| 10 | CMR | Cameroon | SSA | WCA | Lower middle income (LM) | 0.07 |
| 17 | CIV | C™te d'Ivoire | SSA | WCA | Lower middle income (LM) | 0.03 |
| 20 | DJI | Diibouti | SSA | ESA | Lower middle income (LM) | 0.09 |

```
In [70]:    1  #the next few graphs will consist of a similar pattern as I just want to
            2  #visualize them in as many ways possible
            3  region_order = ['EAP', 'ECA', 'LAC', 'MENA', 'SA', 'SSA']
            4
            5  plt.figure(figsize=(12, 6))
            6  ax1 = sns.boxplot(x='Region', y='Total', data=dftest2task3, order=region_order)
            7  plt.title('Box Plot of Secondary Children with Internet Access for Lower Middle incon
            8  plt.xlabel('Region')
            9  plt.ylabel('Total')
           10
           11  ax1.set_ylim(0, 1)
           12
           13  plt.show()
           14
           15  plt.figure(figsize=(12, 6))
           16  ax2 = sns.boxplot(x='Region', y='Total', data=dftesttask3, order=region_order)
           17  plt.title('Box Plot of Primary Children with Internet Access for Lower Middle income
           18  plt.xlabel('Region')
           19  plt.ylabel('Total')
           20
           21  ax2.set_ylim(0, 1)
           22
           23  plt.show()
           24
```

Box Plot of Secondary Children with Internet Access for Lower Middle income countries by Region

Box Plot of Primary Children with Internet Access for Lower Middle income countries by Region

```
In [71]:  1  region_order = ['EAP', 'ECA', 'LAC', 'MENA', 'SA', 'SSA']
          2
          3  plt.figure(figsize=(12, 6))
          4  ax1 = sns.barplot(x='Region', y='Total', data=dftest2task3, order=region_order)
          5  plt.title('Bar Plot of Secondary Children with Internet Access for Lower Middle incor
          6  plt.xlabel('Region')
          7  plt.ylabel('Total')
          8
          9  ax1.set_ylim(0, 1)
         10
         11  plt.show()
         12
         13  plt.figure(figsize=(12, 6))
         14  ax2 = sns.barplot(x='Region', y='Total', data=dftesttask3, order=region_order)
         15  plt.title('Bar Plot of Primary Children with Internet Access for Lower Middle income
         16  plt.xlabel('Region')
         17  plt.ylabel('Total')
         18
         19  ax2.set_ylim(0, 1)
         20
         21  plt.show()
         22
```
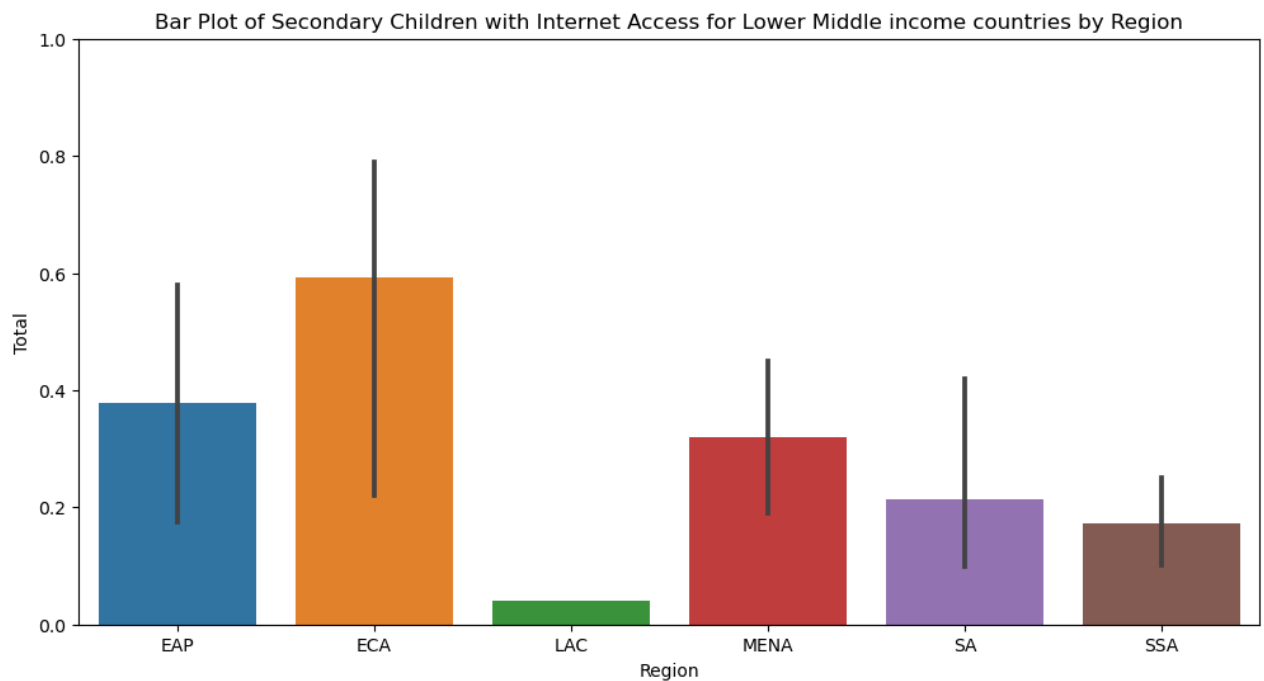


Bar Plot of Secondary Children with Internet Access for Lower Middle income countries by Region

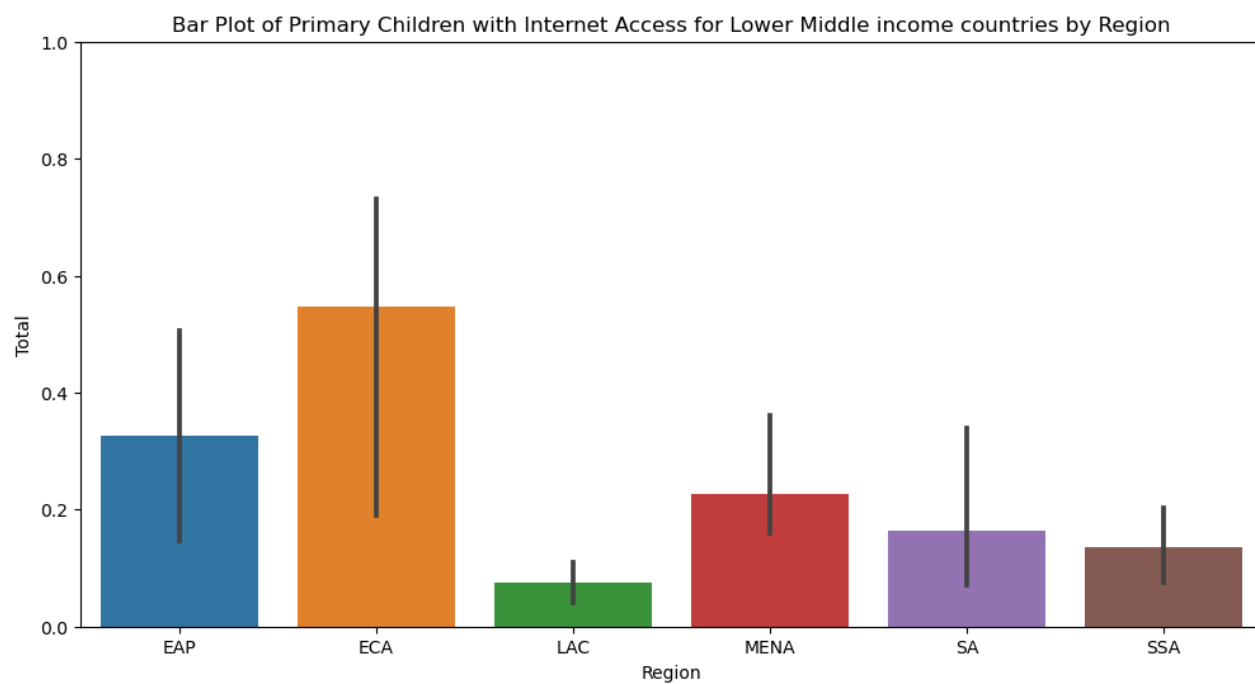Bar Plot of Primary Children with Internet Access for Lower Middle income countries by Region

```
 1  region_order = ['EAP', 'ECA', 'LAC', 'MENA', 'SA', 'SSA']
 2
 3  plt.figure(figsize=(12, 6))
 4  ax1 = sns.violinplot(x='Region', y='Total', data=dftest2task3, order=region_order)
 5  plt.title('Violin Plot of Secondary Children with Internet Access for Lower Middle i
 6  plt.xlabel('Region')
 7  plt.ylabel('Total')
 8
 9  ax1.set_ylim(-0.4, 1.5)
10
11  plt.show()
12
13  plt.figure(figsize=(12, 6))
14  ax2 = sns.violinplot(x='Region', y='Total', data=dftesttask3, order=region_order)
15  plt.title('Violin Plot of Primary Children with Internet Access for Lower Middle inc
16  plt.xlabel('Region')
17  plt.ylabel('Total')
18
19  ax2.set_ylim(-0.4, 1.5)
20
21  plt.show()
22
```



Violin Plot of Secondary Children with Internet Access for Lower Middle income countries by Region

Violin Plot of Primary Children with Internet Access for Lower Middle income countries by Region

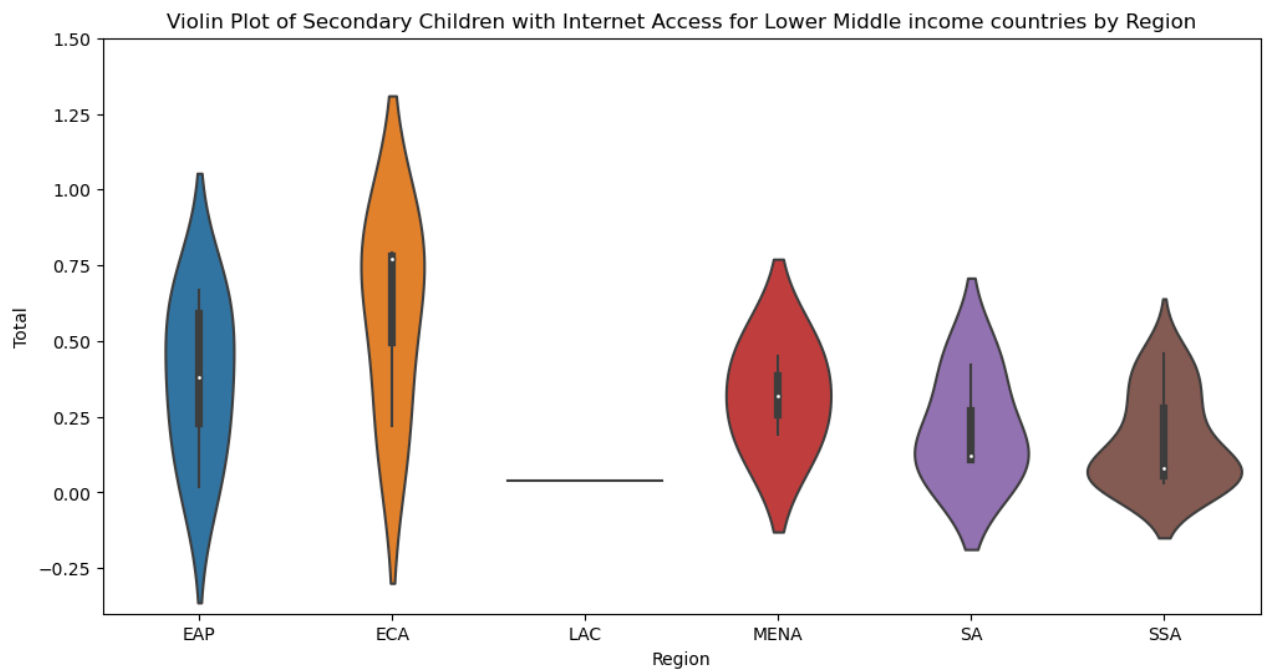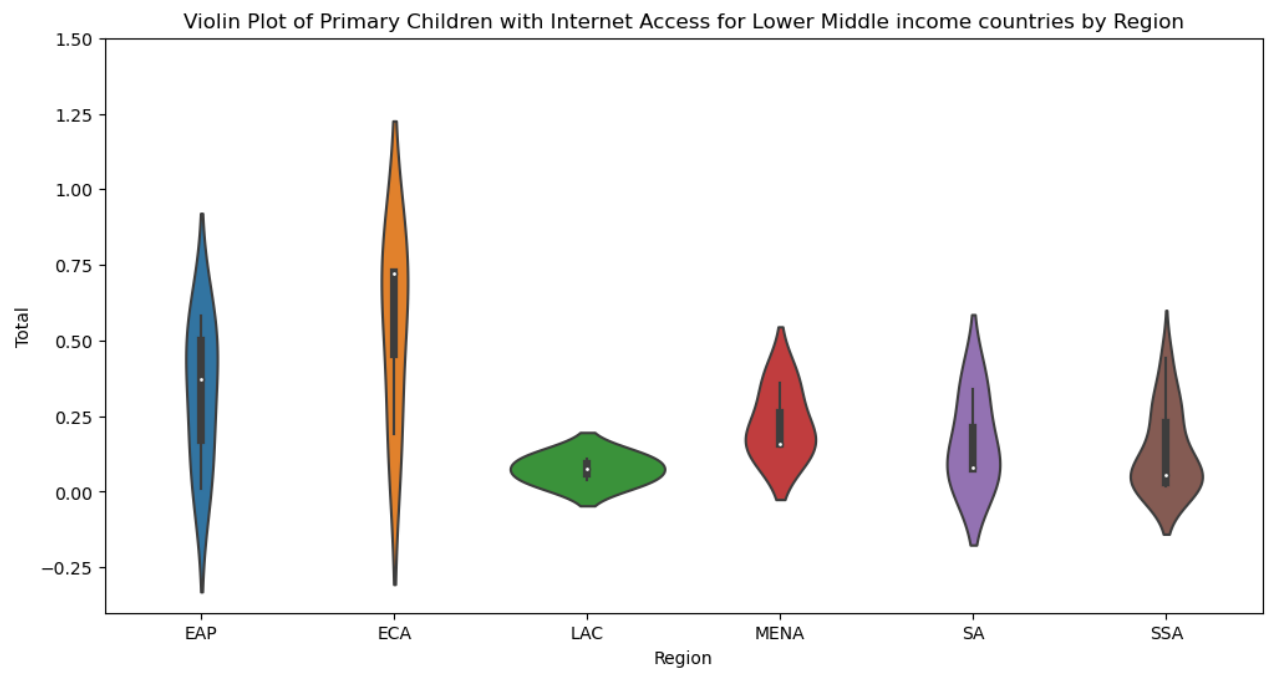```
In [73]:   1
           2  region_order = ['EAP', 'ECA', 'LAC', 'MENA', 'SA', 'SSA']
           3
           4  plt.figure(figsize=(12, 6))
           5  ax1 = sns.swarmplot(x='Region', y='Total', data=dftest2task3, order=region_order)
           6  plt.title('Swarm Plot of Secondary Children with Internet Access for Lower Middle inc
           7  plt.xlabel('Region')
           8  plt.ylabel('Total')
           9
          10  ax1.set_ylim(0, 1)
          11
          12  plt.show()
          13
          14  plt.figure(figsize=(12, 6))
          15  ax2 = sns.swarmplot(x='Region', y='Total', data=dftesttask3, order=region_order)
          16  plt.title('Swarm Plot of Primary Children with Internet Access for Lower Middle incor
          17  plt.xlabel('Region')
          18  plt.ylabel('Total')
          19
          20  ax2.set_ylim(0, 1)
          21
          22  plt.show()
          23
```



Swarm Plot of Secondary Children with Internet Access for Lower Middle income countries by Region

Swarm Plot of Primary Children with Internet Access for Lower Middle income countries by Region