Submitted by: Angelo Luis C. Cu

```
!pip install hvplot
```

```
Requirement already satisfied: hvplot in /usr/local/lib/python3.10/dist-packages (0.9.2)
Requirement already satisfied: bokeh>=1.0.0 in /usr/local/lib/python3.10/dist-packages (from hvplot) (3.3.4)
Requirement already satisfied: colorcet>=2 in /usr/local/lib/python3.10/dist-packages (from hvplot) (3.1.0)
Requirement already satisfied: holoviews>=1.11.0 in /usr/local/lib/python3.10/dist-packages (from hvplot) (1.17.1)
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from hvplot) (2.0.3)
Requirement already satisfied: numpy>=1.15 in /usr/local/lib/python3.10/dist-packages (from hvplot) (1.25.2)
Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from hvplot) (24.0)
Requirement already satisfied: panel>=0.11.0 in /usr/local/lib/python3.10/dist-packages (from hvplot) (1.3.8)
Requirement already satisfied: param<3.0,>=1.12.0 in /usr/local/lib/python3.10/dist-packages (from hvplot) (2.1.0)
Requirement already satisfied: Jinja2>=2.9 in /usr/local/lib/python3.10/dist-packages (from bokeh>=1.0.0->hvplot) (3.1.3)
Requirement already satisfied: contourpy>=1 in /usr/local/lib/python3.10/dist-packages (from bokeh>=1.0.0->hvplot) (1.2.1)
Requirement already satisfied: pillow>=7.1.0 in /usr/local/lib/python3.10/dist-packages (from bokeh>=1.0.0->hvplot) (9.4.0)
Requirement already satisfied: PyYAML>=3.10 in /usr/local/lib/python3.10/dist-packages (from bokeh>=1.0.0->hvplot) (6.0.1)
Requirement already satisfied: tornado>=5.1 in /usr/local/lib/python3.10/dist-packages (from bokeh>=1.0.0->hvplot) (6.3.3)
Requirement already satisfied: xyzservices>=2021.09.1 in /usr/local/lib/python3.10/dist-packages (from bokeh>=1.0.0->hvplot) (2024.4.0)
Requirement already satisfied: pyviz-comms>=0.7.4 in /usr/local/lib/python3.10/dist-packages (from holoviews>=1.11.0->hvplot) (3.0.2)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas->hvplot) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas->hvplot) (2023.4)
Requirement already satisfied: tzdata>=2022.1 in /usr/local/lib/python3.10/dist-packages (from pandas->hvplot) (2024.1)
Requirement already satisfied: markdown in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (3.6)
Requirement already satisfied: markdown-it-py in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (3.0.0)
Requirement already satisfied: linkify-it-py in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (2.0.3)
Requirement already satisfied: mdit-py-plugins in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (0.4.0)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (2.31.0)
Requirement already satisfied: tqdm>=4.48.0 in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (4.66.2)
Requirement already satisfied: bleach in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (6.1.0)
Requirement already satisfied: typing-extensions in /usr/local/lib/python3.10/dist-packages (from panel>=0.11.0->hvplot) (4.11.0)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from Jinja2>=2.9->bokeh>=1.0.0->hvplot) (2.1.
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas->hvplot) (1.16.0
Requirement already satisfied: webencodings in /usr/local/lib/python3.10/dist-packages (from bleach->panel>=0.11.0->hvplot) (0.5.1)
Requirement already satisfied: uc-micro-py in /usr/local/lib/python3.10/dist-packages (from linkify-it-py->panel>=0.11.0->hvplot) (1.0.3
Requirement already satisfied: mdurl~=0.1 in /usr/local/lib/python3.10/dist-packages (from markdown-it-py->panel>=0.11.0->hvplot) (0.1.2
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->panel>=0.11.0->hvplot
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->panel>=0.11.0->hvplot) (3.7)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->panel>=0.11.0->hvplot) (2.0
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->panel>=0.11.0->hvplot) (202
```

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

from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.linear_model import LinearRegression

%matplotlib inline
```

```python
life_df = pd.read_csv('/content/data/Life Expectancy Data.csv')
life_df
```

| | Country | Year | Status | Life expectancy | Adult Mortality | infant deaths | Alcohol | percentage expenditure | Hepatitis B | Measles | ... | Polio | Total expenditure | D |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 2015 | Developing | 65.0 | 263.0 | 62 | 0.01 | 71.279624 | 65.0 | 1154 | ... | 6.0 | 8.16 | |
| 1 | Afghanistan | 2014 | Developing | 59.9 | 271.0 | 64 | 0.01 | 73.523582 | 62.0 | 492 | ... | 58.0 | 8.18 | |
| 2 | Afghanistan | 2013 | Developing | 59.9 | 268.0 | 66 | 0.01 | 73.219243 | 64.0 | 430 | ... | 62.0 | 8.13 | |
| 3 | Afghanistan | 2012 | Developing | 59.5 | 272.0 | 69 | 0.01 | 78.184215 | 67.0 | 2787 | ... | 67.0 | 8.52 | |
| 4 | Afghanistan | 2011 | Developing | 59.2 | 275.0 | 71 | 0.01 | 7.097109 | 68.0 | 3013 | ... | 68.0 | 7.87 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 2933 | Zimbabwe | 2004 | Developing | 44.3 | 723.0 | 27 | 4.36 | 0.000000 | 68.0 | 31 | ... | 67.0 | 7.13 | |
| 2934 | Zimbabwe | 2003 | Developing | 44.5 | 715.0 | 26 | 4.06 | 0.000000 | 7.0 | 998 | ... | 7.0 | 6.52 | |
| 2935 | Zimbabwe | 2002 | Developing | 44.8 | 73.0 | 25 | 4.43 | 0.000000 | 73.0 | 304 | ... | 73.0 | 6.53 | |
| 2936 | Zimbabwe | 2001 | Developing | 45.3 | 686.0 | 25 | 1.72 | 0.000000 | 76.0 | 529 | ... | 76.0 | 6.16 | |
| 2937 | Zimbabwe | 2000 | Developing | 46.0 | 665.0 | 24 | 1.68 | 0.000000 | 79.0 | 1483 | ... | 78.0 | 7.10 | |

2938 rows × 22 columns

## Data Wrangling

```python
# checks for duplicate values
life_df[life_df.duplicated()].shape[0]
```

```
0
```

```python
# changes spaces to underscores for easier column access
life_df.columns = [column.replace(' ', '_') for column in life_df.columns]
life_df.columns = [column.strip('_') for column in life_df.columns]
life_df
```

| | Country | Year | Status | Life_expectancy | Adult_Mortality | infant_deaths | Alcohol | percentage_expenditure | Hepatitis_B | Measles |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 2015 | Developing | 65.0 | 263.0 | 62 | 0.01 | 71.279624 | 65.0 | 1154 |
| 1 | Afghanistan | 2014 | Developing | 59.9 | 271.0 | 64 | 0.01 | 73.523582 | 62.0 | 492 |
| 2 | Afghanistan | 2013 | Developing | 59.9 | 268.0 | 66 | 0.01 | 73.219243 | 64.0 | 430 |
| 3 | Afghanistan | 2012 | Developing | 59.5 | 272.0 | 69 | 0.01 | 78.184215 | 67.0 | 2787 |
| 4 | Afghanistan | 2011 | Developing | 59.2 | 275.0 | 71 | 0.01 | 7.097109 | 68.0 | 3013 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 2933 | Zimbabwe | 2004 | Developing | 44.3 | 723.0 | 27 | 4.36 | 0.000000 | 68.0 | 31 |
| 2934 | Zimbabwe | 2003 | Developing | 44.5 | 715.0 | 26 | 4.06 | 0.000000 | 7.0 | 998 |
| 2935 | Zimbabwe | 2002 | Developing | 44.8 | 73.0 | 25 | 4.43 | 0.000000 | 73.0 | 304 |
| 2936 | Zimbabwe | 2001 | Developing | 45.3 | 686.0 | 25 | 1.72 | 0.000000 | 76.0 | 529 |
| 2937 | Zimbabwe | 2000 | Developing | 46.0 | 665.0 | 24 | 1.68 | 0.000000 | 79.0 | 1483 |

2938 rows × 22 columns

```
# checks for missing values
life_df.info()
# life expectancy, adult mortality, alcohol, hepatitis B, BMI,
# polio-diptheria, gdp-schooling have NaN values
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2938 entries, 0 to 2937
Data columns (total 22 columns):
 #   Column                           Non-Null Count  Dtype
---  ------                           --------------  -----
 0   Country                          2938 non-null   object
 1   Year                             2938 non-null   int64
 2   Status                           2938 non-null   object
 3   Life_expectancy                  2928 non-null   float64
 4   Adult_Mortality                  2928 non-null   float64
 5   infant_deaths                    2938 non-null   int64
 6   Alcohol                          2744 non-null   float64
 7   percentage_expenditure           2938 non-null   float64
 8   Hepatitis_B                      2385 non-null   float64
 9   Measles                          2938 non-null   int64
 10  BMI                              2904 non-null   float64
 11  under-five_deaths                2938 non-null   int64
 12  Polio                            2919 non-null   float64
 13  Total_expenditure                2712 non-null   float64
 14  Diphtheria                       2919 non-null   float64
 15  HIV/AIDS                         2938 non-null   float64
 16  GDP                              2490 non-null   float64
 17  Population                       2286 non-null   float64
 18  thinness__1-19_years             2904 non-null   float64
 19  thinness_5-9_years               2904 non-null   float64
 20  Income_composition_of_resources  2771 non-null   float64
 21  Schooling                        2775 non-null   float64
dtypes: float64(16), int64(4), object(2)
memory usage: 505.1+ KB
```

```
# as all of the missing values are from columns with numerical values,
# and the missing values are less than half of the total count,
# I decided to fill them with their mean
# since country is available, I decided to group them by country to try to minimize bias

for column in life_df.columns:
  if life_df[column].dtype != 'object':
    life_df[column] = life_df.groupby('Country')[column].transform(lambda x: x.fillna(x.mean()))

life_df
```

| | Country | Year | Status | Life_expectancy | Adult_Mortality | infant_deaths | Alcohol | percentage_expenditure | Hepatitis_B | Measles |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Afghanistan | 2015 | Developing | 65.0 | 263.0 | 62 | 0.01 | 71.279624 | 65.0 | 1154 |
| 1 | Afghanistan | 2014 | Developing | 59.9 | 271.0 | 64 | 0.01 | 73.523582 | 62.0 | 492 |
| 2 | Afghanistan | 2013 | Developing | 59.9 | 268.0 | 66 | 0.01 | 73.219243 | 64.0 | 430 |
| 3 | Afghanistan | 2012 | Developing | 59.5 | 272.0 | 69 | 0.01 | 78.184215 | 67.0 | 2787 |
| 4 | Afghanistan | 2011 | Developing | 59.2 | 275.0 | 71 | 0.01 | 7.097109 | 68.0 | 3013 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 2933 | Zimbabwe | 2004 | Developing | 44.3 | 723.0 | 27 | 4.36 | 0.000000 | 68.0 | 31 |
| 2934 | Zimbabwe | 2003 | Developing | 44.5 | 715.0 | 26 | 4.06 | 0.000000 | 7.0 | 998 |
| 2935 | Zimbabwe | 2002 | Developing | 44.8 | 73.0 | 25 | 4.43 | 0.000000 | 73.0 | 304 |
| 2936 | Zimbabwe | 2001 | Developing | 45.3 | 686.0 | 25 | 1.72 | 0.000000 | 76.0 | 529 |
| 2937 | Zimbabwe | 2000 | Developing | 46.0 | 665.0 | 24 | 1.68 | 0.000000 | 79.0 | 1483 |

2938 rows × 22 columns

```
life_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2938 entries, 0 to 2937
Data columns (total 22 columns):
 #   Column                           Non-Null Count  Dtype
```

```
 ---  ------                                 --------------  -----
  0   Country                                2938 non-null   object
  1   Year                                   2938 non-null   int64
  2   Status                                 2938 non-null   object
  3   Life_expectancy                        2928 non-null   float64
  4   Adult_Mortality                        2928 non-null   float64
  5   infant_deaths                          2938 non-null   int64
  6   Alcohol                                2921 non-null   float64
  7   percentage_expenditure                 2938 non-null   float64
  8   Hepatitis_B                            2794 non-null   float64
  9   Measles                                2938 non-null   int64
  10  BMI                                    2904 non-null   float64
  11  under-five_deaths                      2938 non-null   int64
  12  Polio                                  2938 non-null   float64
  13  Total_expenditure                      2906 non-null   float64
  14  Diphtheria                             2938 non-null   float64
  15  HIV/AIDS                               2938 non-null   float64
  16  GDP                                    2533 non-null   float64
  17  Population                             2290 non-null   float64
  18  thinness__1-19_years                   2904 non-null   float64
  19  thinness_5-9_years                     2904 non-null   float64
  20  Income_composition_of_resources        2771 non-null   float64
  21  Schooling                              2775 non-null   float64
 dtypes: float64(16), int64(4), object(2)
 memory usage: 505.1+ KB
```

```
# since there are countries without data for that specific column at all,
# I decided to fill them with the general mean

for column in life_df.columns:
  if life_df[column].dtype != 'object':
    life_df[column] = life_df[column].fillna(life_df[column].mean())

life_df.info()
# all missing data are handled
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2938 entries, 0 to 2937
Data columns (total 22 columns):
 #   Column                                 Non-Null Count  Dtype
 ---  ------                                 --------------  -----
  0   Country                                2938 non-null   object
  1   Year                                   2938 non-null   int64
  2   Status                                 2938 non-null   object
  3   Life_expectancy                        2938 non-null   float64
  4   Adult_Mortality                        2938 non-null   float64
  5   infant_deaths                          2938 non-null   int64
  6   Alcohol                                2938 non-null   float64
  7   percentage_expenditure                 2938 non-null   float64
  8   Hepatitis_B                            2938 non-null   float64
  9   Measles                                2938 non-null   int64
  10  BMI                                    2938 non-null   float64
  11  under-five_deaths                      2938 non-null   int64
  12  Polio                                  2938 non-null   float64
  13  Total_expenditure                      2938 non-null   float64
  14  Diphtheria                             2938 non-null   float64
  15  HIV/AIDS                               2938 non-null   float64
  16  GDP                                    2938 non-null   float64
  17  Population                             2938 non-null   float64
  18  thinness__1-19_years                   2938 non-null   float64
  19  thinness_5-9_years                     2938 non-null   float64
  20  Income_composition_of_resources        2938 non-null   float64
  21  Schooling                              2938 non-null   float64
 dtypes: float64(16), int64(4), object(2)
 memory usage: 505.1+ KB
```

```
# converting categorical to numerical data
columns = [
    'Country','Status'
] # columns to get the unique values
unique_values = []

# gets the unique values of a column and appends it to the unique_values list
for column in columns:
  unique_values.append(life_df[column].unique().tolist())
unique_values
```

```
[['Afghanistan',
  'Albania',
  'Algeria',
  'Angola',
  'Antigua and Barbuda',
  'Argentina',
  'Armenia',
  'Australia',
  'Austria',
  'Azerbaijan',
  'Bahamas',
  'Bahrain',
  'Bangladesh',
  'Barbados',
  'Belarus',
  'Belgium',
  'Belize',
  'Benin',
  'Bhutan',
  'Bolivia (Plurinational State of)',
  'Bosnia and Herzegovina',
  'Botswana',
  'Brazil',
  'Brunei Darussalam',
  'Bulgaria',
  'Burkina Faso',
  'Burundi',
  "Côte d'Ivoire",
  'Cabo Verde',
```

```
        'Cambodia',
        'Cameroon',
        'Canada',
        'Central African Republic',
        'Chad',
        'Chile',
        'China',
        'Colombia',
        'Comoros',
        'Congo',
        'Cook Islands',
        'Costa Rica',
        'Croatia',
        'Cuba',
        'Cyprus',
        'Czechia',
        "Democratic People's Republic of Korea",
        'Democratic Republic of the Congo',
        'Denmark',
        'Djibouti',
        'Dominica',
        'Dominican Republic',
        'Ecuador',
        'Egypt',
        'El Salvador',
        'Equatorial Guinea',
        'Eritrea',
        'Estonia',
        'Ethiopia',
```

```python
# creates the dictionaries
result_dicts = [] # stores the results here

for data in unique_values:
  keys = [i for i in data]
  values = [i for i in range(1, len(data)+1)]
  result_dicts.append({keys[i] : values[i] for i in range(len(values))})
result_dicts
```

```
        'Qatar': 137,
        'Republic of Korea': 138,
        'Republic of Moldova': 139,
        'Romania': 140,
        'Russian Federation': 141,
        'Rwanda': 142,
        'Saint Kitts and Nevis': 143,
        'Saint Lucia': 144,
        'Saint Vincent and the Grenadines': 145,
        'Samoa': 146,
        'San Marino': 147,
        'Sao Tome and Principe': 148,
        'Saudi Arabia': 149,
        'Senegal': 150,
        'Serbia': 151,
        'Seychelles': 152,
        'Sierra Leone': 153,
        'Singapore': 154,
        'Slovakia': 155,
        'Slovenia': 156,
        'Solomon Islands': 157,
        'Somalia': 158,
        'South Africa': 159,
        'South Sudan': 160,
        'Spain': 161,
        'Sri Lanka': 162,
        'Sudan': 163,
        'Suriname': 164,
        'Swaziland': 165,
        'Sweden': 166,
        'Switzerland': 167,
        'Syrian Arab Republic': 168,
        'Tajikistan': 169,
        'Thailand': 170,
        'The former Yugoslav republic of Macedonia': 171,
        'Timor-Leste': 172,
        'Togo': 173,
        'Tonga': 174,
        'Trinidad and Tobago': 175,
        'Tunisia': 176,
        'Turkey': 177,
        'Turkmenistan': 178,
        'Tuvalu': 179,
        'Uganda': 180,
        'Ukraine': 181,
        'United Arab Emirates': 182,
        'United Kingdom of Great Britain and Northern Ireland': 183,
        'United Republic of Tanzania': 184,
        'United States of America': 185,
        'Uruguay': 186,
        'Uzbekistan': 187,
        'Vanuatu': 188,
        'Venezuela (Bolivarian Republic of)': 189,
        'Viet Nam': 190,
        'Yemen': 191,
        'Zambia': 192,
        'Zimbabwe': 193},
      {'Developing': 1, 'Developed': 2}]
```

```python
# maps the categorical data to their numerical counterparts
for column in range(len(columns)):
  life_df.replace(result_dicts[column], inplace=True)

life_df
```

|  | Country | Year | Status | Life_expectancy | Adult_Mortality | infant_deaths | Alcohol | percentage_expenditure | Hepatitis_B | Measles | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2015 | 1 | 65.0 | 263.0 | 62 | 0.01 | 71.279624 | 65.0 | 1154 | ... |
| 1 | 1 | 2014 | 1 | 59.9 | 271.0 | 64 | 0.01 | 73.523582 | 62.0 | 492 | ... |
| 2 | 1 | 2013 | 1 | 59.9 | 268.0 | 66 | 0.01 | 73.219243 | 64.0 | 430 | ... |
| 3 | 1 | 2012 | 1 | 59.5 | 272.0 | 69 | 0.01 | 78.184215 | 67.0 | 2787 | ... |
| 4 | 1 | 2011 | 1 | 59.2 | 275.0 | 71 | 0.01 | 7.097109 | 68.0 | 3013 | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 2933 | 193 | 2004 | 1 | 44.3 | 723.0 | 27 | 4.36 | 0.000000 | 68.0 | 31 | ... |
| 2934 | 193 | 2003 | 1 | 44.5 | 715.0 | 26 | 4.06 | 0.000000 | 7.0 | 998 | ... |
| 2935 | 193 | 2002 | 1 | 44.8 | 73.0 | 25 | 4.43 | 0.000000 | 73.0 | 304 | ... |
| 2936 | 193 | 2001 | 1 | 45.3 | 686.0 | 25 | 1.72 | 0.000000 | 76.0 | 529 | ... |
| 2937 | 193 | 2000 | 1 | 46.0 | 665.0 | 24 | 1.68 | 0.000000 | 79.0 | 1483 | ... |

2938 rows × 22 columns

## Exploratory Data Analysis

```
life_df.describe()
"""
According to this data, the average country is developing,
with a life expectancy of 69 years and a population of 12.73 Million
"""
```

|  | Country | Year | Status | Life_expectancy | Adult_Mortality | infant_deaths | Alcohol | percentage_expenditure | Hepatii |
|---|---|---|---|---|---|---|---|---|---|
| count | 2938.000000 | 2938.000000 | 2938.000000 | 2938.000000 | 2938.000000 | 2938.000000 | 2938.000000 | 2938.000000 | 2938.00 |
| mean | 96.091219 | 2007.518720 | 1.174268 | 69.224932 | 164.796448 | 30.303948 | 4.600849 | 738.251295 | 78.64 |
| std | 56.250042 | 4.613841 | 0.379405 | 9.507640 | 124.080302 | 117.926501 | 4.027279 | 1987.914858 | 24.55 |
| min | 1.000000 | 2000.000000 | 1.000000 | 36.300000 | 1.000000 | 0.000000 | 0.010000 | 0.000000 | 1.00 |
| 25% | 47.000000 | 2004.000000 | 1.000000 | 63.200000 | 74.000000 | 0.000000 | 0.930000 | 4.685343 | 73.50 |
| 50% | 94.000000 | 2008.000000 | 1.000000 | 72.000000 | 144.000000 | 3.000000 | 3.780000 | 64.912906 | 88.00 |
| 75% | 146.000000 | 2012.000000 | 1.000000 | 75.600000 | 227.000000 | 22.000000 | 7.677500 | 441.534144 | 96.00 |
| max | 193.000000 | 2015.000000 | 2.000000 | 89.000000 | 723.000000 | 1800.000000 | 17.870000 | 19479.911610 | 99.00 |

8 rows × 22 columns

```
# comparing by developed and developing countries
developing_life = life_df.query('Status == 1')
developed_life = life_df.query('Status == 2')


developed_life.mean()

    Country                          9.525309e+01
    Year                             2.007523e+03
    Status                           1.000000e+00
    Life_expectancy                  6.712018e+01
    Adult_Mortality                  1.827588e+02
    infant_deaths                    3.638417e+01
    Alcohol                          3.493100e+00
    percentage_expenditure           3.234703e+02
    Hepatitis_B                      7.745355e+01
    Measles                          2.824926e+03
    BMI                              3.547577e+01
    under-five_deaths                5.052514e+01
    Polio                            8.000298e+01
    Total_expenditure                5.576311e+00
    Diphtheria                       7.980067e+01
    HIV/AIDS                         2.088664e+00
    GDP                              4.668433e+03
    Population                       1.374722e+07
    thinness__1-19_years             5.582378e+00
    thinness_5-9_years               5.624522e+00
    Income_composition_of_resources  5.845291e-01
    Schooling                        1.125592e+01
    dtype: float64


developing_life.mean()

    Country                          1.000625e+02
    Year                             2.007500e+03
    Status                           2.000000e+00
    Life_expectancy                  7.919785e+01
    Adult_Mortality                  7.968555e+01
    infant_deaths                    1.494141e+00
    Alcohol                          9.849678e+00
    percentage_expenditure           2.703600e+03
    Hepatitis_B                      8.430827e+01
    Measles                          4.990059e+02
    BMI                              5.180391e+01
    under-five_deaths                1.810547e+00
    Polio                            9.373633e+01
    Total_expenditure                7.554042e+00
    Diphtheria                       9.347656e+01
```

```
HIV/AIDS                          1.000000e-01
GDP                               2.021901e+04
Population                        7.937177e+06
thinness__1-19_years              1.320703e+00
thinness_5-9_years                1.296680e+00
Income_composition_of_resources   8.314013e-01
Schooling                         1.548429e+01
dtype: float64
```

```python
plt.figure(figsize=(20,20))
sns.heatmap(
    life_df.sort_index().corr(),
    annot=True, center=0, square=True
)
```