

# Quantifying Wellbeing: A Data Science Approach to Understanding Human Longevity

2025-06-20

## Project Motivation

This project was conducted independently to analyze socio-economic factors affecting Life Expectancy at Birth (LEB). Using a real-world dataset, the aim was to explore linear and non-linear relationships and build predictive models. The project uses Python and statistical techniques such as Pearson and Spearman correlations, linear regression, and data transformation to draw conclusions and validate model accuracy.

## Step 1: Splitting the Data

To begin the analysis, I split the dataset into two subsets using the `train_test_split` function from scikit-learn. The training set comprises **75%** of the data, while the remaining **25%** is reserved for testing. A random seed (`SEED = 123456`) ensures reproducibility. All columns, including `Country`, were retained to allow flexible post-analysis.

## Step 2: Single-Variable Model

### a) Identifying the Strongest Correlate to LEB

Using the Pearson correlation coefficient, I identified the Human Development Index (HDI) as the variable with the strongest linear relationship to LEB.

Table 1: Variable with the Strongest Linear Relationship with LEB

Variable	Pearson Correlation Coefficient
Human Development Index (value)	0.93

### b) Building the Linear Regression Model

A simple linear regression model was created using HDI as the predictor and LEB as the target. The model was trained on the training set and yielded the following results:

Table 2: Linear Model Results

Metric	Value
Slope	48.85
Intercept	36.86
$R^2$ (Coefficient of Determination)	0.87

Linear Regression: Life Expectancy at Birth, both sexes (years) vs Human Development Index (value)

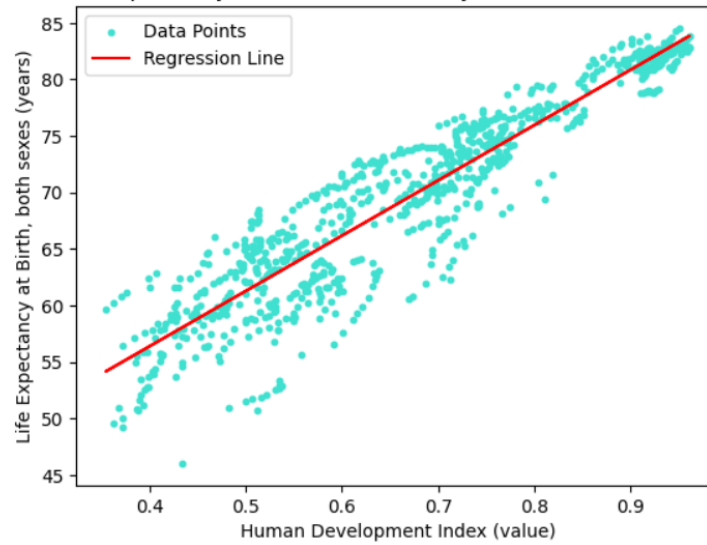


Figure 1: Scatter plot showing relationship between HDI and LEB

### c) Test Set Prediction and Evaluation

The model was applied to the test set to evaluate its performance. I calculated both the Pearson correlation between predicted and actual LEB values and the Mean Squared Error (MSE).

Table 3: Model Performance Metrics on the Test Set

Metric	Value
Correlation between predicted and actual LEB	0.9133
Mean Squared Error (MSE)	11.4083

### d) Interpretation

Since HDI incorporates life expectancy, education, and income levels, its strong relationship with LEB is intuitive. These elements are known to influence health outcomes. Higher HDI is generally associated with improved living standards and healthcare access, which directly affects life expectancy [2, 1].

### Step 3: Exploring Non-Linear Relationships

To identify non-linear associations, I visualized scatter plots of all variables against LEB. One such relationship was between Gross National Income (GNI) and LEB, which suggested a logarithmic curve.

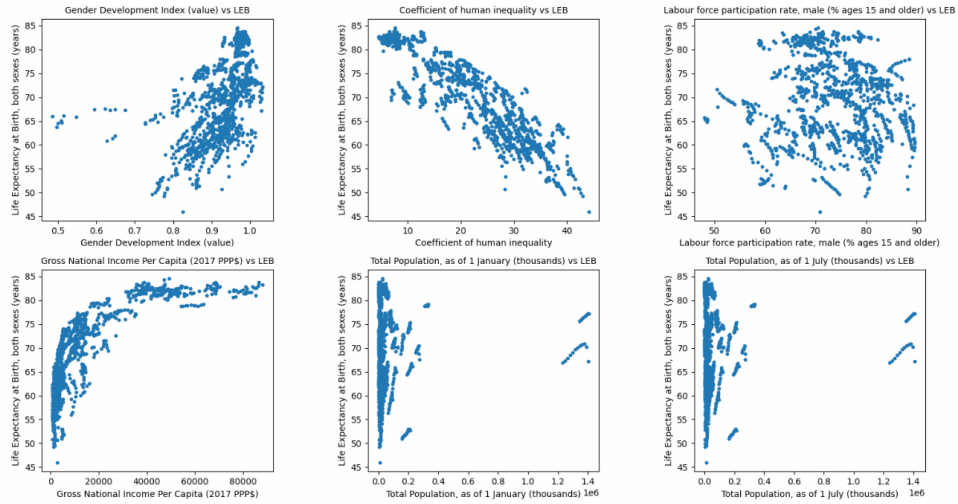


Figure 2: Scatter plots between LEB and selected variables

The correlations for GNI before and after logarithmic transformation are as follows:

Table 4: Correlation between GNI and LEB (Before Transformation)

Correlation Type	Correlation Value
Pearson Correlation	0.77
Spearman Correlation	0.89

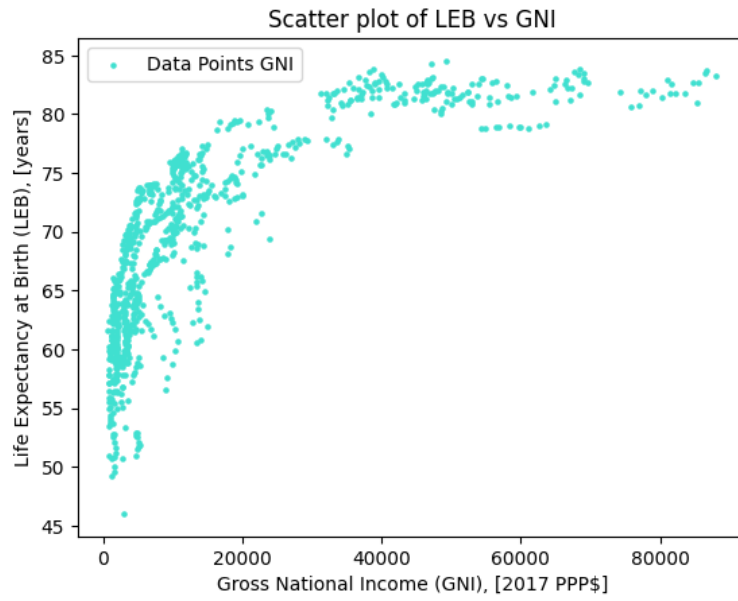


Figure 3: LEB vs GNI (before log transformation)

Table 5: Correlation after Logarithmic Transformation

Correlation Type	Correlation Value
Pearson Correlation (after transformation)	0.88

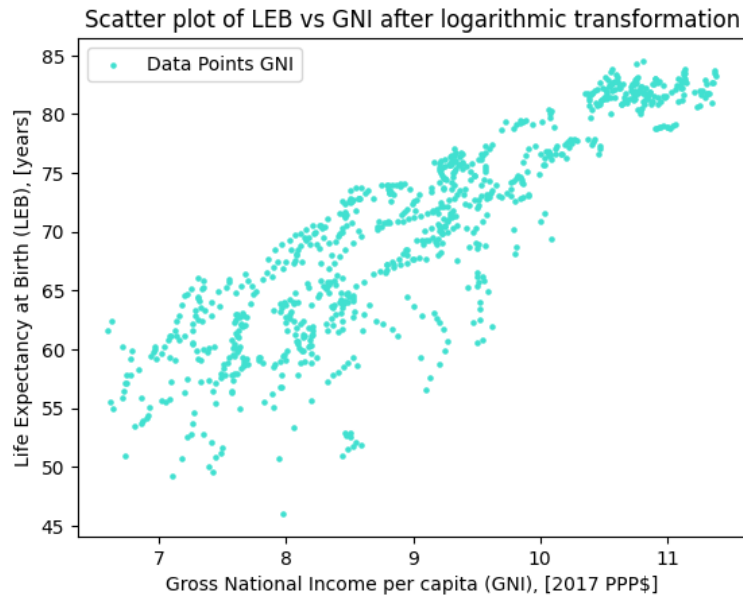


Figure 4: LEB vs  $\log(\text{GNI})$

## Step 4: Multiple Linear Regression Model

To enhance prediction performance, I constructed a multiple linear regression model using the top three most correlated features with LEB:

- Median Age
- Expected Years of Schooling, female
- Expected Years of Schooling (overall)

Table 6: Top Variables Correlated with LEB

Variable	Pearson Correlation
Median Age	0.8955
Expected Years of Schooling, female	0.8453
Expected Years of Schooling (overall)	0.8261

Table 7: Multiple Linear Model Coefficients

Metric	Value
Coefficients	0.72, 2.64, -2.25
Intercept	46.84
R <sup>2</sup> (Training)	0.836

Table 8: Model Performance Metrics on Test Set

Metric	Value
Correlation between predicted and actual LEB	0.9209
Mean Squared Error (MSE)	10.2688

## Conclusion

This project demonstrated that both single and multiple regression models can predict life expectancy reasonably well using socio-economic indicators. The HDI was the strongest single predictor, but using a combination of education and demographic variables led to slightly improved performance. Transforming variables like GNI further enhanced correlation, validating the non-linear nature of some relationships.

## References

- UNDP: Human Development Reports. <https://hdr.undp.org>
- Our World in Data: Life Expectancy. <https://ourworldindata.org/life-expectancy>

## References

- [1] Our World in Data. *Life Expectancy: How is it calculated and how should it be interpreted?* Retrieved 2024-11-29. URL: <https://ourworldindata.org/life-expectancy-how-is-it-calculated-and-how-should-it-be-interpreted>.
- [2] UNDP (United Nations Development Programme). *Human Development Index (HDI) Data Center*. Retrieved 2024-11-27. 2023. URL: <https://hdr.undp.org/data-center/human-development-index#/indicies/HDI>.

## Appendix

### Problem 1 - Importing necessary libraries and splitting the data

```
1 #Import necessary libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LinearRegression
7 from sklearn.metrics import mean_squared_error
8 from scipy.stats import pearsonr, spearmanr
9 from sklearn.metrics import r2_score
10
11 #Connect to drive to be able to import files
12 from google.colab import drive
13 drive.mount('/content/drive')
14
15 #Load the csv file and dropping empty rows
16 df = pd.read_csv('/content/drive/MyDrive/
17                 life_expectancy.csv').dropna()
18 print(df.columns)
19
20 #Split the dataset into training and test
21 SEED = 123456
22 df_train, df_test = train_test_split(df, random_state=
23                                     SEED, test_size=0.25)
24 print("Training data size:", df_train.shape)
25 print("Testing data size:", df_test.shape)
26
27 Dataset = len(df)
28 TrainingSet = len(df_train)
29 TestSet= len(df_test)
30
31 print(f"Dataset (Length: {Dataset}) divided to \
32       nTraining Set (Length: {TrainingSet}) and Test Set \
33       (Length: {TestSet}) ")
```

### Problem 2 - Single-variable model

#### a) Variable with the strongest linear relationship to LEB

```
1 #Target variable
2 LEB = "Life Expectancy at Birth, both sexes (years)"
3
4 #only numeric columns are selected
5 df_train = df_train.select_dtypes(include='number')
6 df_train = df_train.dropna()
7
8 #Calculate the correlation for all variables with LEB
9 pearson = df_train.corr(method='pearson')
10
```

```

11 #Drop the correlation of LEB with itself
12 pearson_drop = pearson[LEB].drop(LEB)
13
14 #Find the maximum correlation and the corresponding
    variable
15 strongest_corr = pearson_drop.max()
16 strongest_var = pearson_drop.loc[pearson_drop ==
    strongest_corr].index[0]
17
18 #Print the strongest variable and its correlation
19 print(f"The variable with the strongest linear
    relationship with LEB is '{strongest_var}' "
20       f"with a Pearson correlation coefficient of {
    strongest_corr:.2f}.")

```

#### b) Single-variable linear model

```

1 #Initialize the Linear Regression model
2 lr = LinearRegression()
3
4 #Define the independent and dependent variables
5 independent_variables = ['Human Development Index (
    value)']
6 dependent_variable = LEB
7
8 #Prepare the training set
9 X_train = df_train[independent_variables]
10 y_train = df_train[dependent_variable]
11
12 #Fit the model to the training data
13 lr.fit(X_train,y_train)
14
15 #Coefficient of determination slope and intercept
16 slope = lr.coef_[0]
17 intercept = lr.intercept_
18
19 # R2 coefficient of determination
20 r2=lr.score(X_train,y_train)
21
22 # calculate the regression line by slope and intercept
23 regression_line = intercept + slope * X_train
24
25 print(f"Slope: {slope:.2f}")
26 print(f"Intercept: {intercept:.2f}")
27 print(f"R2 (Coefficient of Determination): {r2:.2f}")
28
29 #Scatter plot showing the LEB as the function of the
    variable
30 plt.scatter(X_train, y_train, color='Turquoise', label
    ='Data Points',s=10)

```



```

31 #plot the regression line
32 plt.plot(X_train, regression_line, color='red', label=
    'Regression Line')
33
34 #Plot the regression line over the scatter plot
35 plt.xlabel('Human Development Index (value)')
36 plt.ylabel('Life Expectancy at Birth, both sexes (
    years)')
37 plt.title(f"Linear Regression: {dependent_variable} vs
    Human Development Index (value)")
38 plt.legend()
39 plt.show()

```

### c) Predict the LEB for the test set

```

1 #Prepare the test set
2 X_test = df_test[independent_variables]
3 y_test = df_test[dependent_variable]
4
5 #predict the LEB for test set
6 y_pred_test = lr.predict(X_test)
7
8 #calculate the correlation between the predicted
    values and the target variable LEB
9 correlation, _ = pearsonr(y_test, y_pred_test)
10
11 # Mean Squared Error (MSE) between the predictions and
    the true target variable values
12 mse = mean_squared_error(y_test, y_pred_test)
13
14 print(f"Correlation between predicted and the target
    variable LEB: {correlation:.4f}")
15 print(f"Mean squared error (MSE): {mse:.4f}")

```

### Problem 3 - Non-linear relationship

```

1 #Plotting all variables compared to LEB for visual
    study
2
3 #Plotting done on all data, not just training data
4
5 dependent_variable = LEB
6 excluded_columns = [LEB, 'Human Development Index (
    value)'] #Removing LEB and HDI
7 independent_variables = [col for col in df.columns if
    col not in excluded_columns]
8
9 #Create subplots
10 n_cols = 3
11 n_rows = 16

```

```

12 fig, axes = plt.subplots(n_rows, n_cols, figsize=(15,
13                               64))
14 #Flatten axes for easy iteration
15 axes = axes.ravel()
16
17 #Plot scatter plots
18 for i, column in enumerate(independent_variables):
19     axes[i].scatter(df[column], df[dependent_variable], s=10)
20     axes[i].set_title(f"{column} vs LEB", fontsize =
21                       10)
22     axes[i].set_xlabel(column, fontsize = 10)
23     axes[i].set_ylabel(dependent_variable, fontsize =
24                       10)
25
26 #Adjust layout and show
27 plt.tight_layout()
28 plt.show()
29
30 #Pearson and spearman coefficient: GNI and LEB
31
32 #Define the independent and dependent variables
33 GNI = "Gross National Income Per Capita (2017 PPP$)"
34 #Using training data
35 independent_variable = df_train[GNI]
36 dependent_variable = df_train[LEB]
37
38 #Calculate Pearson and Spearman correlations
39 pearson_corr, _ = pearsonr(independent_variable,
40                             dependent_variable)
41 spearman_corr, _ = spearmanr(independent_variable,
42                              dependent_variable)
43
44 #Print results
45 print(f"Correlation between GNI and LEB:")
46 print(f"    Pearson Correlation: {pearson_corr:.2f}")
47 print(f"    Spearman Correlation: {spearman_corr:.2f}")
48
49 #Scatter plot showing the LEB as the function of the
50 GNI
51 plt.scatter(independent_variable, dependent_variable,
52             color='Turquoise', label='Data Points GNI',s=5)
53
54 #Plot
55 plt.xlabel('Gross National Income (GNI), [2017 PPP$]')
56 plt.ylabel('Life Expectancy at Birth (LEB), [years]')
57 plt.title(f"Scatter plot of LEB vs GNI")
58 plt.legend()
59 plt.show()

```

```

54 #Logarithmic transformation on GNI data
55 independent_variable_transformed = np.log(
    independent_variable)
56
57 #Calculate Pearson and Spearman correlations after
    transformation
58 pearson_corr_transformed, _ = pearsonr(
    independent_variable_transformed,
    dependent_variable)
59
60 #Print results
61 print(f"Correlation between GNI and LEB after
    transformation:")
62 print(f"    Pearson Correlation: {
    pearson_corr_transformed:.2f}")
63
64 #Scatter plot showing the LEB as the function of the
    GNI after transformation
65 plt.scatter(independent_variable_transformed,
    dependent_variable, color='Turquoise', label='Data
    Points GNI',s=5)
66
67 #Plot
68 plt.xlabel('Gross National Income per capita (GNI),
    [2017 PPP$]')
69 plt.ylabel('Life Expectancy at Birth (LEB), [years]')
70 plt.title(f"Scatter plot of LEB vs GNI after
    logarithmic transformation")
71 plt.legend()
72 plt.show()

```

#### Problem 4 - Multiple linear regression

```

1 #Finding variables with high pearson correlation to
    LEB
2
3 #Only numeric columns are selected
4 df_train = df_train.select_dtypes(include='number')
5 df_train = df_train.dropna()
6
7 #Calculate the correlation for all variables with LEB
8 pearson = df_train.corr(method='pearson')
9
10 #Drop the correlation of LEB with itself and HDI
11 pearson_drop = pearson[LEB].drop([LEB, 'Human
    Development Index (value)'])
12
13 #Sort in descending order
14 sorted_pearson = pearson_drop.sort_values(ascending=
    False)

```

```

15 print(sorted_pearson)
16
17
18 #Initialize the Linear Regression model
19 lr = LinearRegression()
20
21 #Define the independent and dependent variables
22 independent_variables = ['Median Age, as of 1 July (
    years)', 'Expected Years of Schooling, female (
    years)', 'Expected Years of Schooling (years)']
23 dependent_variable = LEB
24
25 #Prepare the training set
26 X_train = df_train[independent_variables]
27 y_train = df_train[dependent_variable]
28
29 #Prepare the test set
30 X_test = df_test[independent_variables]
31 y_test = df_test[dependent_variable]
32
33 #Fit the model to the training data
34 lr.fit(X_train,y_train)
35
36 #Coefficients intercept
37 coefficients = lr.coef_
38 intercept = lr.intercept_
39
40 #R^2 coefficient of determination
41 r2=lr.score(X_train,y_train)
42
43 print("Coefficients:", coefficients)
44 print(f"Intercept: {intercept:.2f}")
45 print(f"R^2 (Coefficient of Determination): {r2:.3f}")
46
47
48 #Predict the LEB for test set
49 y_pred_test = lr.predict(X_test)
50
51 #calculate the correlation between the predicted
    values and the target variable LEB
52 correlation, _ = pearsonr(y_test, y_pred_test)
53
54 # Mean Squared Error (MSE) between the predictions and
    the true target variable values
55 mse = mean_squared_error(y_test, y_pred_test)
56
57 print(f"Pearson Correlation between predicted and the
    target variable LEB: {correlation:.4f}")
58 print(f"Mean squared error (MSE): {mse:.4f}")

```