

# DAT565/DIT407 Assignment 4

David Duong  
thuan@chalmers.se

Savinjith Walisadeera  
savwal@chalmers.se

September 30, 2024

## Problem 1

We perform a 75:25 train-test split randomly of the data. We use this ratio because it gives us a sufficient amount of points to train the model and test it. We select randomly to make sure we do not fit specifically to any trends.

## Problem 2

- a) We identified the strongest linear relationship to LEB was Human Development Index. This parameter had the highest Pearson product-moment correlation coefficient.
- b) The coefficient of determination was 0.8483

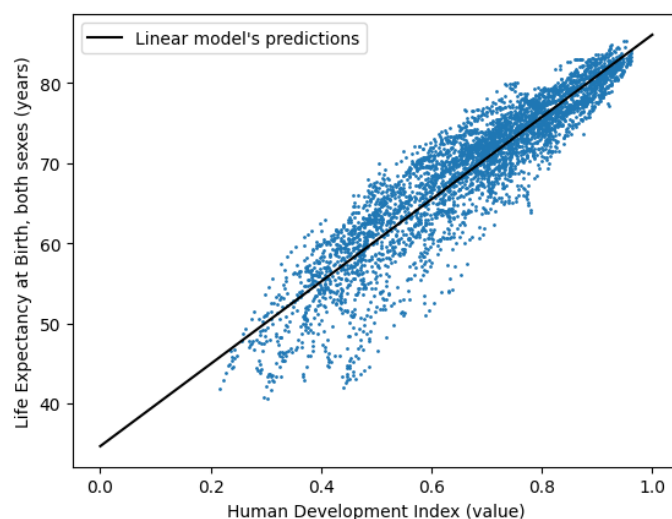


Figure 1: Scatter plot of human development index to life expectancy after birth with a single-variable linear model.

- c) The correlation between HDI and LEB is: 0.912. The mean square error between the the test data and trained model was 14.60.

d) Human development index is a coefficient which measures: healthy and long lives, education level and standard of living, which are normalized to compute the index between 0 and 1 [1]. These factors are play a large role for the general public of a countries life expectancy after birth. Therefore it makes sense that HDI has a strong correlation coefficient with LEB.

### Problem 3

We look at large difference in size between the min and max values and found that GDP per capita and LEB seemed logarithmically correlated. Pearson coefficient before transformation: 0.651. Pearson coefficient after transformation: 0.837.

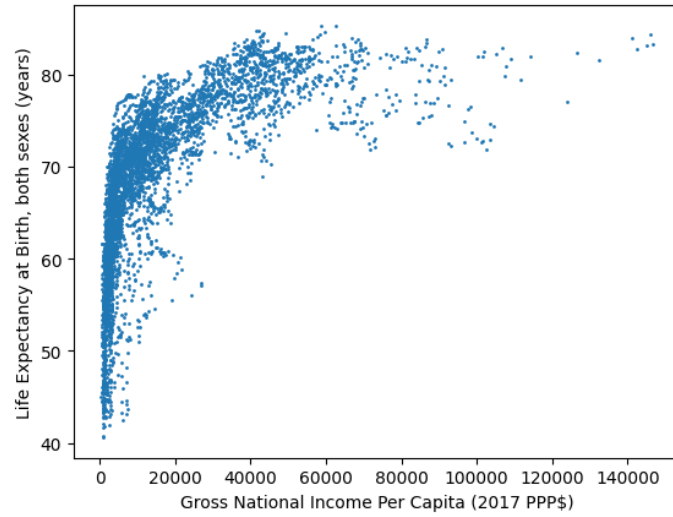


Figure 2: Scatter plot of GDP per capita to life expectancy after birth before transformation.

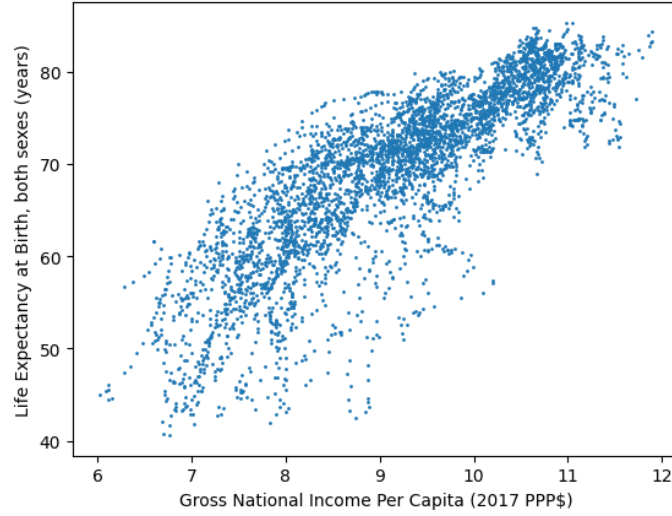


Figure 3: Scatter plot of GDP per capita to life expectancy after birth after transformation.

## Problem 4

We chose the variables by using our knowledge on what impacts a country's life expectancy after birth combined with the absolute value of the Pearson product-moment correlation coefficient between all parameters against life expectancy after birth. With these variables we created a multiple linear regression model. The variables in the model are: expected years of schooling, crude birth rate, net reproduction rate gross national income per capita and gender development index. When we had missing values for a variable we omitted that row from the data set.

We have the Pearson correlation = 0.887,  $R^2=0.867$ , mean square error = 11.821, and the intercept value = 65.864.

Variables	k
Expected Years of Schooling	0.492
Crude birth rate	-1.198
Net reproduction rate	13.932
Gross national income per capita	1.873
Gender development index	-11.234

Table 1: Table of coefficients for each variable in the model.

This model had better model evaluation parameters than the single variable linear regression.

## References

- [1] Human Development Index (HDI) *United Nations Development Program*, 2024.

## A Code

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import r2_score
5 from sklearn.metrics import mean_squared_error
6 import numpy as np
7 import matplotlib.pyplot as plt
8
9 df = pd.read_csv("life_expectancy.csv")
10 df
11
12 #Split the data 75-25
13 train_data, test_data = train_test_split(df, test_size
    =0.25, random_state=12)
14 train_data
15
16 data_no_country = train_data.drop('Country', axis=1)
17 cov_matrix = data_no_country.cov()
18 variances = data_no_country.var()
19 standard_deviations = np.sqrt(variances)
20
21
22 pearson_coef = cov_matrix['Life_Expectancy_at_Birth,
    both_sexes(years)']/(standard_deviations*
    standard_deviations['Life_Expectancy_at_Birth,
    both_sexes(years)'])
23
24 #We find that Human Development Index (value) has the
    strongest linear relationship with LEB
25 pearson_coef.sort_values()
26
27 train_LEB = train_data['Life_Expectancy_at_Birth,
    both_sexes(years)']
28 train_HDI = train_data['Human_Development_Index(value)']
29
30 strongest = train_data[['Life_Expectancy_at_Birth,
    both_sexes(years)', 'Human_Development_Index(value)']]
31
32 train_HDI.values.reshape(-1,1)
```

```

33 model = LinearRegression().fit(train_HDI.values.
    reshape(-1,1), train_LEB.values)
34 print(model.coef_, model.intercept_)
35 R2 = r2_score(train_LEB.values.reshape(-1,1), model.
    predict(train_HDI.values.reshape(-1,1)))
36 R2
37
38 plt.scatter(train_HDI, train_LEB, s=1)
39 plt.plot(np.linspace(start=0, stop=1, num=1000), model.
    predict(np.linspace(start=0, stop=1, num=1000).
    reshape(-1,1)), color='black', label="Linear_models
    'predictions")
40 plt.ylabel('Life Expectancy at Birth, both sexes (
    years)')
41 plt.xlabel('Human Development Index (value)')
42 plt.legend()
43
44 test_LEB = test_data['Life Expectancy at Birth, both
    sexes (years)']
45 test_HDI = test_data['Human Development Index (value)
    ']
46
47 test_predicted_vals = model.predict(test_HDI.values.
    reshape(-1,1))
48 test = pd.DataFrame(test_LEB)
49 test_combined = test.assign(predicted_values=
    test_predicted_vals)
50 test_covariance = test_combined.cov()
51 var_test = test_LEB.var()
52 var_pred = test_predicted_vals.var()
53 test_corr = test_covariance['predicted_values']['Life
    Expectancy at Birth, both sexes (years)']/(np.sqrt(
    var_test*var_pred))
54
55 print('correlation:', test_corr)
56 print('MSE:', mean_squared_error(test_LEB.values,
    test_predicted_vals))
57
58 #Find large differences
59 i = 0
60 vals = []
61 for col in train_data.drop('Country', axis=1):
62     diff = train_data.drop('Country', axis=1).max(axis
        =0)[col]-train_data.drop('Country', axis=1).min
        (axis=0)[col]
63     vals.append((diff, i))
64     i += 1
65
66 def return_first(list):
67     return list[0]

```

```

68 vals.sort(reverse=True,key=return_first)
69 vals
70
71 plt.scatter(np.log(train_data['Gross_National_Income_
    Per_Capita_(2017_PPP$)']),train_data['Life_
    Expectancy_at_Birth,_both_sexes_(years)'],s=1)
72 plt.xlabel('Gross_National_Income_Per_Capita_(2017_
    PPP$)')
73 plt.ylabel('Life_Expectancy_at_Birth,_both_sexes_(
    years)')
74
75 print('Pearson_coefficient_before_transformation:',
    pearson_coef['Gross_National_Income_Per_Capita_
    (2017_PPP$)'])
76
77 log_GDP = np.log(train_data['Gross_National_Income_Per
    _Capita_(2017_PPP$)'])
78
79 dframe = pd.DataFrame(train_LEB)
80 combined = dframe.assign(GDPPC=log_GDP)
81 GDP_covariance = combined.cov()
82 var_train = train_LEB.var()
83 var_GDP = log_GDP.var()
84
85 test_corr = GDP_covariance['GDPPC']['Life_Expectancy_
    at_Birth,_both_sexes_(years)']/(np.sqrt(var_train*
    var_GDP))
86
87 print('Pearson_coefficient_after_transformation:',
    test_corr)
88
89 #Expected Years of Schooling (years)
90 #Gross National Income Per Capita (2017 PPP$)
91 #Crude Birth Rate (births per 1,000 population)
92 #Adolescent Birth Rate (births per 1,000 women ages
    15-19)
93 #Net Reproduction Rate (surviving daughters per woman)
94 #Mean Years of Schooling (years)
95 parameters = ['Expected_Years_of_Schooling_(years)',
96               'Gross_National_Income_Per_Capita_(2017_
    PPP$)',
97               'Crude_Birth_Rate_(births_per_1,000_
    population)',
98               'Net_Reproduction_Rate_(surviving_
    daughters_per_woman)',
99               'Gender_Development_Index_(value)']
100
101 train_data = train_data.dropna(subset=parameters)
102
103

```

```

104 train_LEB = train_data['Life_Expectancy_at_Birth,_both
    _sexes_(years)']
105
106 train_EYS = train_data['Expected_Years_of_Schooling_(
    years)']
107 train_GNIPC = np.log(train_data['Gross_National_Income
    _Per_Capita_(2017_PPP$)'])
108 train_CBR = train_data['Crude_Birth_Rate_(births_per_
    1,000_population)']
109 train_NRR = train_data['Net_Reproduction_Rate_(
    surviving_daughters_per_woman)']
110 train_GDI = train_data['Gender_Development_Index_(
    value)']
111
112 arr = np.array([train_EYS, train_GNIPC, train_CBR,
    train_NRR, train_GDI])
113 arr = np.transpose(arr)
114
115 model = LinearRegression().fit(arr, train_LEB.values)
116 print(model.coef_, model.intercept_)
117
118 R2 = r2_score(train_LEB.values.reshape(-1,1), model.
    predict(arr))
119 print('R2-score: ', R2)
120
121
122 multi_train_predicted_vals = model.predict(arr)
123 multi_train = pd.DataFrame(train_LEB)
124 multi_train_combined = multi_train.assign(
    predicted_values=multi_train_predicted_vals)
125 train_covariance = multi_train_combined.cov()
126 var_train = train_LEB.var()
127 var_pred = multi_train_predicted_vals.var()
128 test_corr = test_covariance['predicted_values']['Life_
    Expectancy_at_Birth,_both_sexes_(years)']/(np.sqrt(
    var_test*var_pred))
129
130 print('correlation: ', test_corr)
131 print('MSE: ', mean_squared_error(train_LEB.values,
    multi_train_predicted_vals))

```