

# DAT565/DIT407 Assignment 4

Saif Sayed  
gussayedfa@student.gu.se

Gona Ibrahim Abdulrahman  
gusibrigo@student.gu.se

2024-05-08

This is a report for assignment 4 for the course *Introduction to Data Science & AI* from Chalmers and Gothenburg University.

## Problem 1: Splitting the data

In problem 1, we utilized the `train_test_split` function provided by the *scikit-learn* library. We start by separating the independent variables, also known as features, from the target variable. In this particular example, the target variable is "Life Expectancy at Birth, both sexes (years)" and we store it in the variable `y`. The remaining columns in the dataset are deemed as features, and we store them in the variable `X`.

To split the data into training and testing sets, we employ the `train_test_split` function. By utilizing this function, we can determine the proportion of the dataset allocated for testing through the `test_size` parameter. For instance, we set the `test_size` to 0.3 implying that 30% of the data will be reserved for testing. To ensure reproducibility and consistency across different runs, we also specify a specific value for the `random_state` parameter, such as 42. This ensures that the same random split is generated each time the code is executed, allowing for reliable comparisons and evaluations.

Our source code can be found in Appendix A of this document.

## Problem 2: Single Variable Model

(A) Here, we calculate the pearson coefficient to identify the variable which has the strongest correlation with the target variable (LEB). The identified variable was 'Human Development Index' with a pearsonr correlation coefficient of 0.92 (2sf).

(B) We developed a linear model for the correlation and found out the coefficient of the model, coefficient of determination and the intercept. And finally plotted a scatter plot with the regression line over it for visualisation (see **Figure 1**). **Table 1** shows our obtained results.

Table 1: Obtained results from our linear model.

Coefficient of determination (2sf)	Coefficient of model (2sf)	Intercept (2sf)
0.84	51.5	34.5

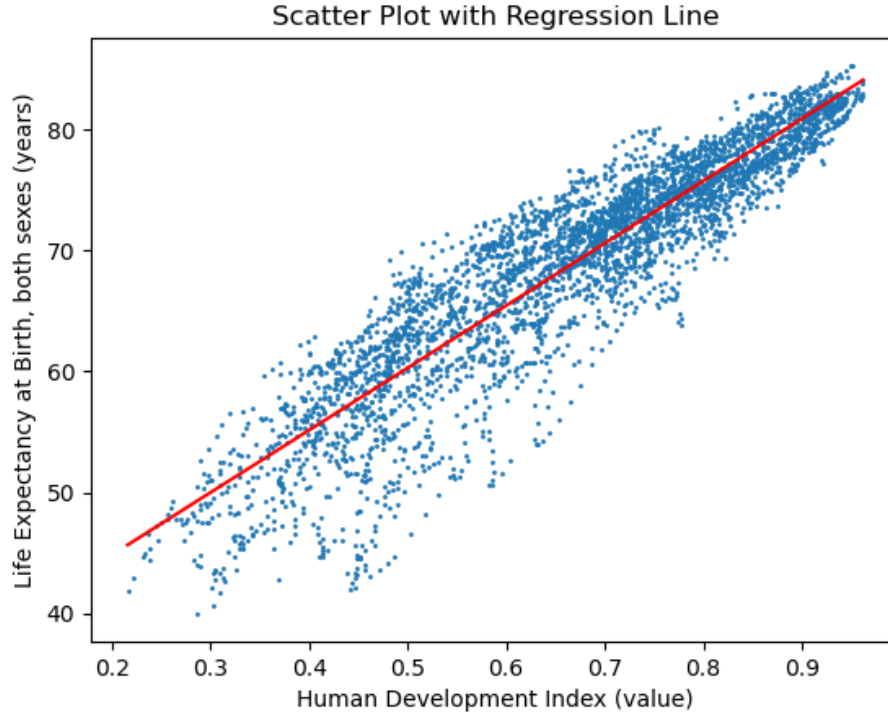


Figure 1: Scatter Plot with regression.

(C) Finally, we predicted the life expectancy for the testing split of our data. Then we calculated the correlation between the predicted and actual values of life expectancy, along with the mean square error. **Table 2** shows our obtained results.

Table 2: Results between predicted and actual LEB.

Pearsonr correlation (2sf)	Mean square error
0.92	12.4

(D) The Human Development Index (HDI) determines the development of countries to ensure the well-being of humans. It takes into consideration a country's economy, education, and healthcare. A higher Human Development Index would mean a better healthcare and education system for its citizens. Thus, humans would have access to quality healthcare that can address health

issues, and education that can raise self-awareness within themselves. A higher economy would also mean a higher income and a better life.

### Problem 3: Non Linear Relationship

In problem 3, we used the spearmanr correlation coefficient to deduce a candidate variable that represents a non-linear relation. We removed the identified variable 'Human Development Index' from the previous task since we already know it exhibits a linear relationship. Unlike the pearsonr correlation coefficient, Spearman's coefficient is designed to capture monotonic relationships, regardless of their linearity. Based on our data analysis, we identified the variable "Gross National Income Per Capita (2017 PPP\$)" as having the highest Spearman's correlation coefficient of 0.87 (2sf). To visualize the correlation between the identified variable and the variable LEB, we created a scatter plot (see **Figure 2**). Upon inspecting the scatter plot, we observed that the relationship between the variables could be represented by a linear pattern after applying a logarithmic transformation to the identified variable (see **Figure 3**). Consequently, we applied the logarithmic function to transform the data. Before the transformation, the Pearson correlation coefficient between the identified variable and LEB was calculated to be 0.65 (2sf). After applying the logarithmic transformation to the identified variable and re-evaluating the correlation, we obtained a Pearson correlation coefficient of 0.83 (2sf).

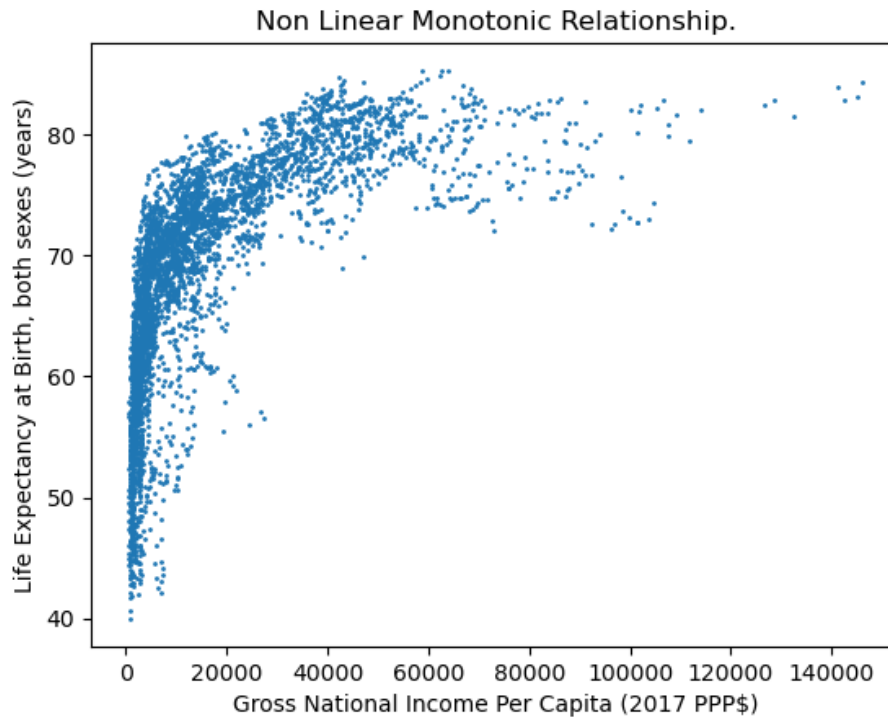


Figure 2: Non Linear Monotonic Relationship.

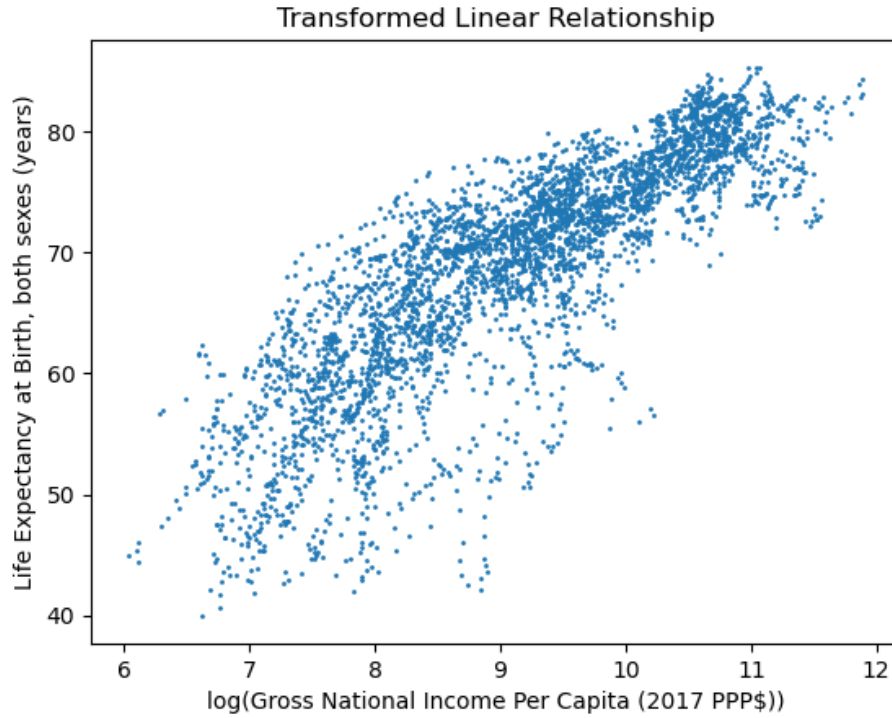


Figure 3: Transformed Linear Relationship.

## Problem 4: Multiple Linear Regression

To identify a subset of variables that can potentially improve the linear model compared to the one in problem 2, we followed a systematic approach. First, we checked the Pearson correlation coefficient between each variable and the target variable. Based on the correlation coefficients, we restricted our selection to the variables exhibiting high correlation ( $> 0.75$ ) with the target variable. These variables have a stronger linear relationship and are more likely to contribute to a better linear model. Finally, we selected a smaller subset of identified variables that are likely to have a meaningful impact on the target variable based on our understanding of the dataset. We calculated the necessary metrics and expanded our set of selected variables if our model didn't outperform the single variable model from problem 2.

The selected variables are:

- Adolescent Birth Rate (births per 1,000 women ages 15-19)
- Expected Years of Schooling (years)
- Median Age, as of 1 July (years)
- Crude Birth Rate (births per 1,000 population)

- Total Fertility Rate (live births per woman)
- Net Reproduction Rate (surviving daughters per woman)

In our multiple linear regression model, we obtained the following results: a coefficient of determination of 0.87 (2sf), an intercept of 58.15 (2sf), and coefficients for the model variables as follows: -0.02, 0.21, 0.40, -0.47, -10.16, 30.52 (2sf). The mean square error of our model was found to be 11.68 (2sf), and the Pearson correlation coefficient between the predicted and actual life expectancy at birth (LEB) was calculated to be 0.93 (in 2sf).

We observed that the mean square error in our multiple linear regression model is slightly lower compared to the mean square error in the single variable model from problem 2. Additionally, we obtained a higher coefficient of determination, indicating a greater proportion of variance explained by our model. These findings suggest that our multiple linear regression model performs slightly better in capturing the relationships between the variables and the target variable.

## Appendix

### A Python code

```

1  import pandas as pd
2  from sklearn.model_selection import train_test_split
3  import matplotlib.pyplot as plt
4  import numpy as np
5  from sklearn.linear_model import LinearRegression
6  from scipy.stats import pearsonr, spearmanr
7  from sklearn.metrics import r2_score
8  from sklearn.metrics import mean_squared_error
9
10 '''
11 Problem 1: Splitting the data
12 '''
13 data = pd.read_csv("life_expectancy.csv", encoding= "
    utf-8")
14
15 # Split the data into features (X) and target variable
    (y)
16 X = data.drop("Life_Expectancy_at_Birth,_both_sexes_(
    years)", axis=1)
17 y = data["Life_Expectancy_at_Birth,_both_sexes_(years)
    "]
18
19 # Perform train-test split
20 X_train, X_test, y_train, y_test = train_test_split(X,
    y, test_size=0.3, random_state=42)
21
22 # Print the shapes of the resulting datasets to verify
    the split

```

```

23 print("Training_set_shape:", X_train.shape, y_train.
      shape)
24 print("Testing_set_shape:", X_test.shape, y_test.shape
      )
25 '''
26 Problem 2: Single Variable Model
27 '''
28 # A.
29
30 # Exclude non-numeric columns from correlation
    calculation
31 numeric_columns = X_train.select_dtypes(include=[float
    , int]).columns
32 numeric_data = X_train[numeric_columns]
33
34 # Handling missing values and replacing them with the
    mean of the corresponding column
35 numeric_data = numeric_data.replace([np.inf, -np.inf],
    np.nan)
36 numeric_data = numeric_data.fillna(numeric_data.mean()
    )
37
38 y_train = y_train.replace([np.inf, -np.inf], np.nan)
39 y_train = y_train.fillna(y_train.mean())
40
41 strongest_variable = ""
42 pearsonr_correlation = 0
43
44 for column in numeric_data.columns:
45     corr= pearsonr(y_train, numeric_data[column])[0]
46     if (corr > pearsonr_correlation):
47         pearsonr_correlation = corr
48         strongest_variable = column
49
50 # Print the results
51 print("Variable_with_the_strongest_correlation:",
    strongest_variable)
52 print("Correlation_coefficient:", pearsonr_correlation
    )
53
54 # B.
55
56 target_variable = "Life_Expectancy_at_Birth,_both_
    sexes_(years)"
57 X = numeric_data[["Human_Development_Index_(value)"]]
58 y = y_train
59
60 # Initialize and fit the linear regression model
61 model = LinearRegression().fit(X, y)
62

```

```

63 actual_LEB = y_train
64 predicted_LEB = model.predict(X)
65
66 determination_coefficient = r2_score(actual_LEB,
    predicted_LEB)
67 model_coefficient = model.coef_
68 intercept = model.intercept_
69
70 print("Coefficient of determination:",
    determination_coefficient)
71 print("Coefficient of model:", model_coefficient)
72 print("Intercept:", intercept)
73
74 # Create a scatter plot
75 plt.scatter(X, y, s=1)
76
77 # Plot the regression line
78 xfit = np.linspace(X["Human_Development_Index_(value)"
    ].min(), X["Human_Development_Index_(value)"].max(),
    , 100)
79 yfit = model.predict(xfit.reshape(-1, 1))
80 plt.plot(xfit, yfit, color='red')
81
82 # Set plot labels
83 plt.xlabel("Human_Development_Index_(value)")
84 plt.ylabel(target_variable)
85 plt.title("Scatter Plot with Regression Line")
86
87 # C.
88
89 test_numeric_columns = X_test.select_dtypes(include=[
    float, int]).columns
90 test_numeric_data = X_test[test_numeric_columns]
91 test_numeric_data = test_numeric_data.replace([np.inf,
    -np.inf], np.nan)
92 test_numeric_data = test_numeric_data.fillna(
    test_numeric_data.mean())
93
94 y_test = y_test.replace([np.inf, -np.inf], np.nan)
95 y_test = y_test.fillna(y_test.mean())
96
97 X_test= test_numeric_data[["Human_Development_Index_(
    value)"]]
98 y_test= y_test
99 X_test = X_test.to_numpy()
100
101 test_y_pred= model.predict(X_test.reshape(-1, 1))
102 pred_target_coefficient = pearsonr(y_test, test_y_pred
    )[0]
103 mean_squared_error= mean_squared_error(y_test,

```

```

        test_y_pred)
104
105 print("The Pearsonr correlation coefficient between
        the predicted LEB and true LEB is:",
        pred_target_coefficient)
106 print("The mean squared error between the predicted
        LEB and true LEB is:", mean_squared_error)
107 '''
108 Problem 3: Non-linear relationship
109 '''
110 numeric_data = numeric_data.drop("Human Development
        Index (value)", axis=1)
111
112 non_linear_variable = ""
113 spearmanr_correlation = 0
114 for column in numeric_data.columns:
115     corr= spearmanr(y_train, numeric_data[column])[0]
116     if (corr > spearmanr_correlation):
117         spearmanr_correlation = corr
118         non_linear_variable = column
119
120 print("The non linear variable is:",
        non_linear_variable)
121 print("Spearmanr Correlation:", spearmanr_correlation)
122
123 x = numeric_data[[non_linear_variable]]
124
125 plt.scatter(x, y, s=1)
126 plt.xlabel(non_linear_variable)
127 plt.ylabel(target_variable)
128 plt.title("Non Linear Monotonic Relationship")
129
130 pearson_coefficient_before = pearsonr(y, x[
        non_linear_variable])[0]
131 print("Pearsonr Correlation Coefficient before:",
        pearson_coefficient_before)
132
133 logarithmic_x = np.log(x)
134
135 plt.scatter(logarithmic_x, y, s=1)
136 plt.xlabel("log(" + non_linear_variable + ")")
137 plt.ylabel(target_variable)
138 plt.title("Transformed Linear Relationship")
139
140 pearson_coefficient_after = pearsonr(y, logarithmic_x[
        non_linear_variable])[0]
141 print("Pearsonr Correlation Coefficient after:",
        pearson_coefficient_after)
142
143 '''

```



```

144 Problem 4: Multiple Linear Regression
145 '''
146 for column in numeric_data.columns:
147     corr= pearsonr(y_train, numeric_data[column])[0]
148     print(column, ":",corr)
149
150 selected_columns= ["Adolescent_Birth_Rate_(births_per_
1,000_women_ages_15-19)", "Expected_Years_of_
Schooling_(years)", "Median_Age,_as_of_1_July_(
years)", "Crude_Birth_Rate_(births_per_1,000_
population)", "Total_Fertility_Rate_(live_births_
per_woman)", "Net_Reproduction_Rate_(surviving_
daughters_per_woman)"]
151 X = numeric_data[selected_columns]
152 y = y_train
153
154 model = LinearRegression().fit(X, y)
155
156 actual_LEB = y_train
157 predicted_LEB = model.predict(X)
158
159 determination_coefficient = r2_score(actual_LEB,
predicted_LEB)
160 model_coefficient = model.coef_
161 intercept = model.intercept_
162
163 print("Coefficient_of_determination:",
determination_coefficient)
164 print("Coefficient_of_model:", model_coefficient)
165 print("Intercept:", intercept)
166
167 X_test= test_numeric_data[selected_columns].values
168 y_test= y_test.values
169
170 test_y_pred= model.predict(X_test.reshape(-1, 6))
171 pred_target_coefficient = pearsonr(y_test, test_y_pred
)[0]
172 print("The_Pearsonr_correlation_coefficient_between_
the_predicted_LEB_and_true_LEB_is:",
pred_target_coefficient)
173
174 '''to run the lines below , the previous lines where
the "mean_squared_error" function was called in
problem 2 needs to be commented out'''
175 #mse = mean_squared_error(y_test, test_y_pred)
176 #print("The mean squared error between the predicted
LEB and true LEB is:", mse)

```