

DIT407 Introduction to data science and AI

Assignment 4

Reynir Siik
reynir.siik@gmail.com

Franco Zambon
guszamfr@student.gu.se

2024-02-19

Abstract

In this assignment we try out different regression models and investigate the underlying data and the quality of the prediction model. When the underlying data for an index is calculated from the index you want to predict a high correlation can be observed. People live longer partly because they live longer. To be careful to use unbiased indexes to base predictions on. Making a non-linear model and taking more indexes into the calculation makes a better model, as shown in problem 3 and 4.

1 Problem 1: Splitting the data

The data is split in two parts to get enough data points in each part to be statistically significant. The splitting is quite straight forward and the code can be easily adapted to splitting in a number of equally large chunks.

2 Problem 2: Single-variable model

2.1 Introduction

After analyzing the Pearson correlation[1] between life expectancy and every other column, the data with the highest correlation, **0.9216**, is Human Development Index. The highest possible correlation (**1.0**) is when the column for life expectancy is correlated to itself, but since that is trivial in this case that result is ignored.

The Human Development Index (HDI) is a summary measure of average achievement in key dimensions of human development: a long and healthy life, being knowledgeable and having a decent standard of living. The HDI is the geometric mean of normalized indices for each of the three dimensions.[2]

The reason for the high correlation is presumed to be due to the underlying data that HDI is calculated from. One factor is Life Expectancy. Another factor is being knowledgeable, which indicates some level of education. *Expected Years of Schooling (years)* has a correlation to *Life Expectancy at Birth* of **0.8202**.

2.2 Training

The prediction graph has Coefficient: **51.06**, and Intercept: **34.83**.

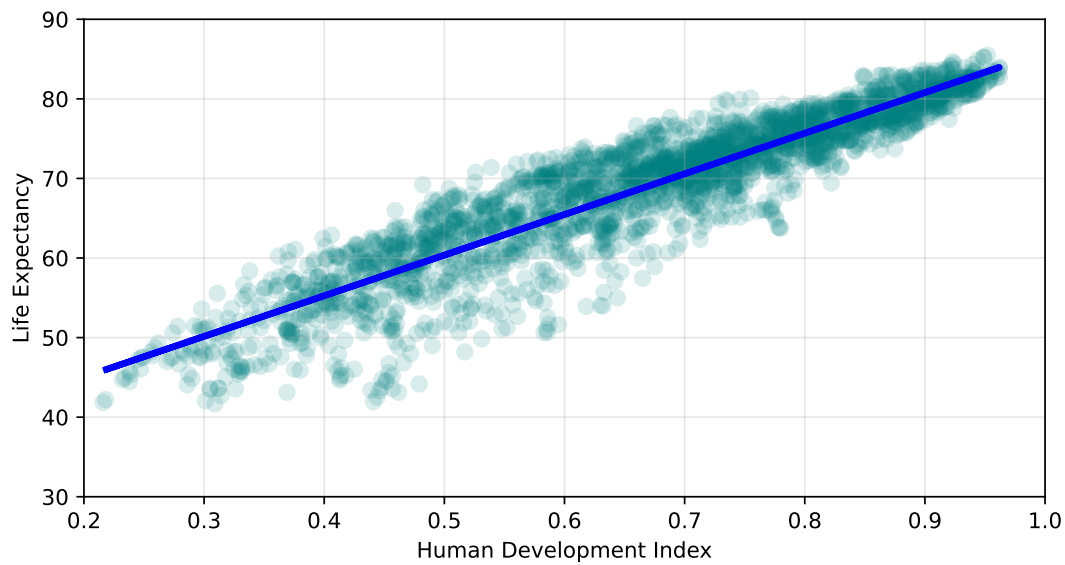


Figure 1: Scatter plot of Human Development Index vs Life Expectancy, training data
Linear prediction model in blue.

2.3 Testing

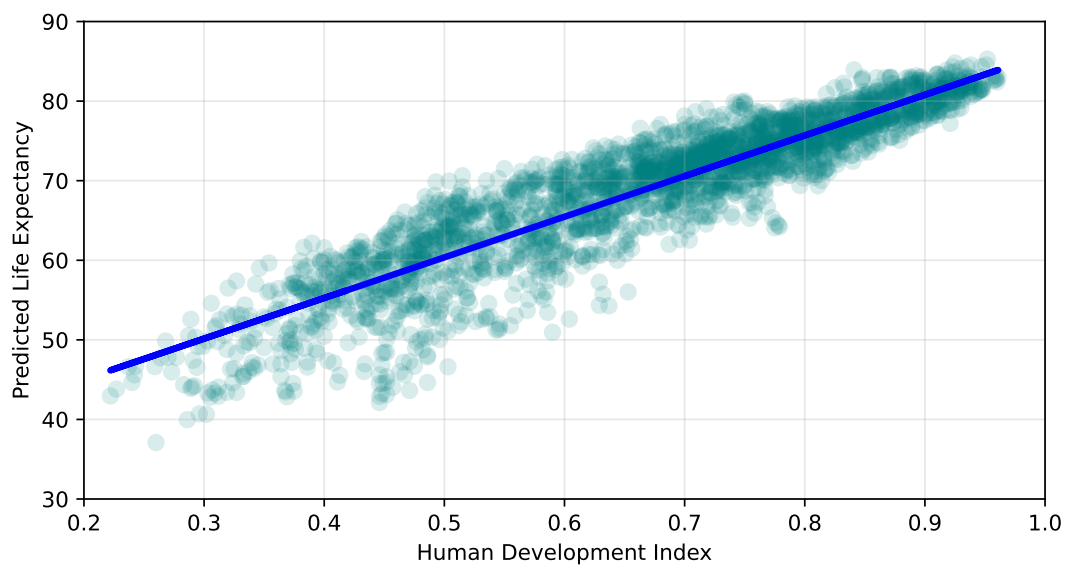


Figure 2: Scatter plot of Human Development Index vs Life Expectancy
Prediction vs real data

The prediction graph versus the test data, figure 2, results in a mean squared error of **13.32**

The coefficient of determination $R^2 = 0.8384$

The correlation between the predicted values and the target variable **0.9157**

3 Problem 3: Non-linear relationship

3.1 Introduction

Looking at the scatter-plots between the different predictors and the dependent variable (Life Expectancy at Birth, both sexes (years)) we found that some of them had a logarithmic monotone relationship with the target variable. We decided to use as a predictor the 'Gross National Income Per Capita (2017 PPP\$)' (Figure 3).

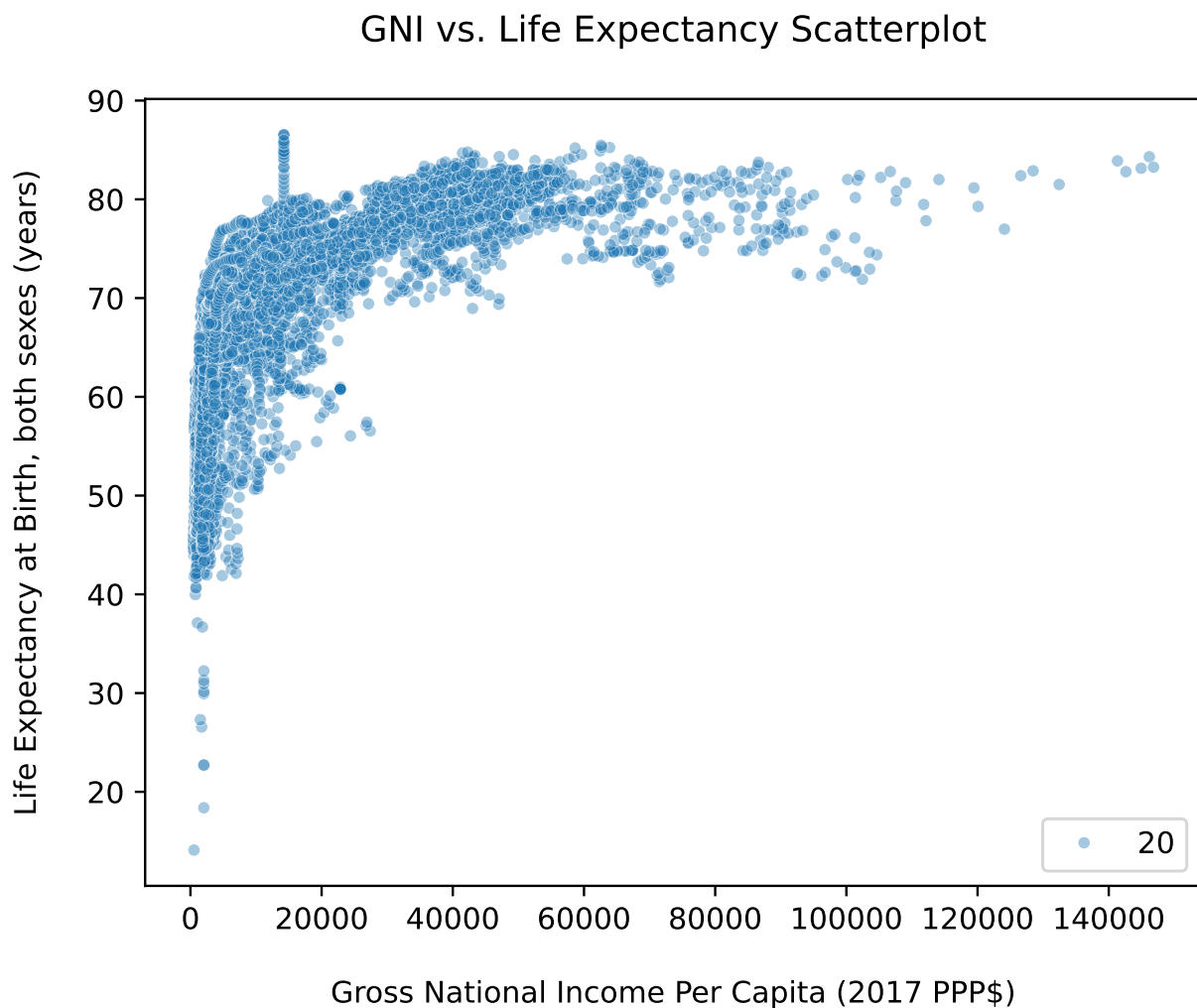


Figure 3: Scatter plot of Gross National Income Per Capita vs Life Expectancy at birth

3.2 Before the transformation

Using this predictor the R-squared that we obtained was 0.4088969210530252, a really low value. Its Pearson coefficient was 0.634651.

3.3 After the transformation

After applying a logarithmic transformation on the x-variable we obtained the scatter plot in figure 4. As we can see the plot shows a much more linear behaviour than before.

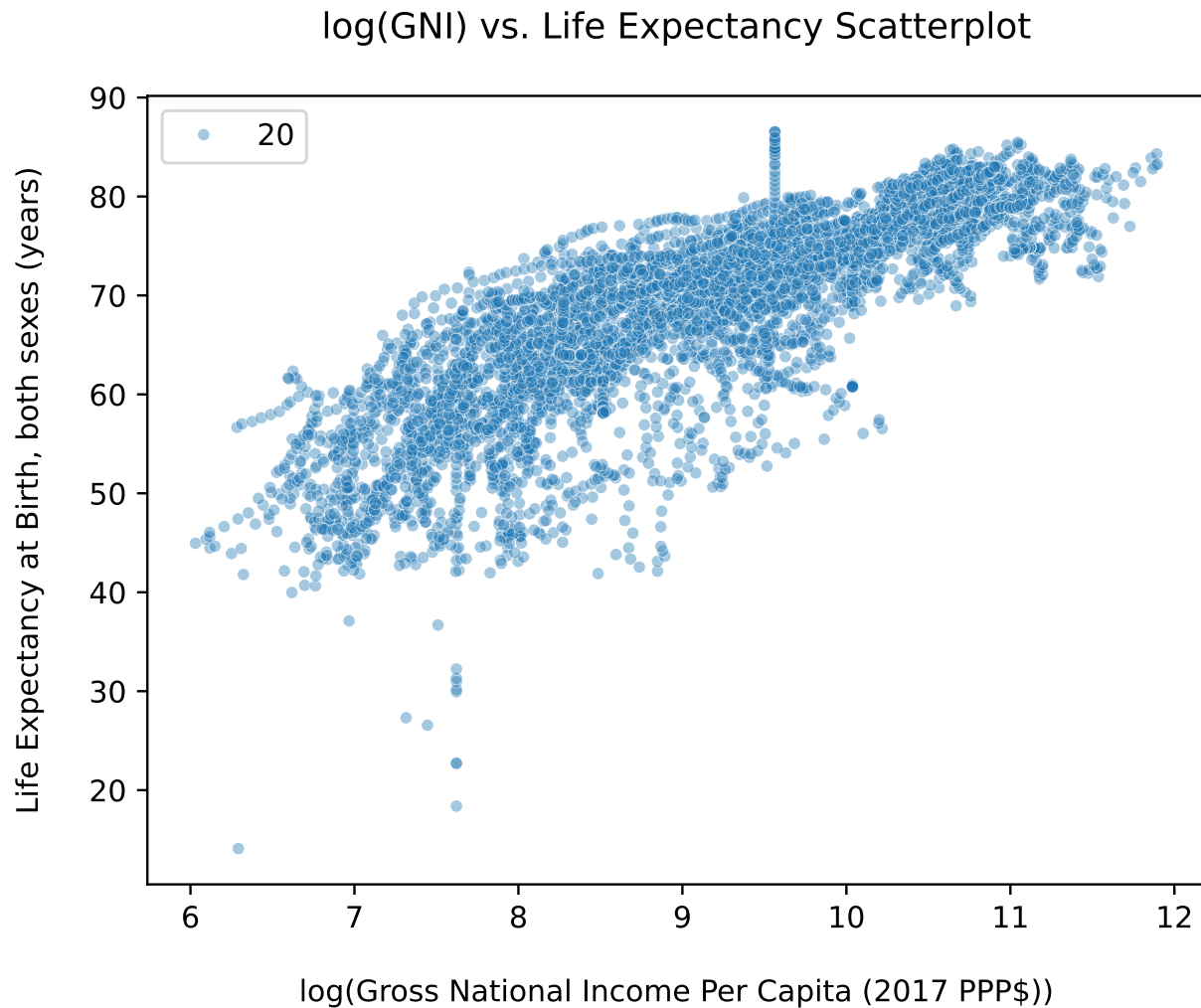


Figure 4: Scatter plot of log(Gross National Income Per Capita) vs Life Expectancy at birth

The R-squared after the transformation is 0.6796613657568554. The Pearson coefficient is 0.8244157723848186.

4 Problem 4: Multiple linear regression

4.1 Introduction

The goal here was to find a minimal subset of variables that allow us to perform a better prediction.

4.2 The process

After printing the list of independent variables and their correlation with the target variable, we selected those which had a coefficient superior to 0.7 or inferior to -0.7 and we tried to apply different combinations of them to the multiple regression model and we evaluated how the the prediction changed just looking at the Adjusted R-squared value.

We also avoided using variables that had a strong correlation between one another (for example we chose to used 'Expected Years of Schooling (years)' but not 'Expected Years of

Schooling, female (years)' nor 'Expected Years of Schooling, male (years)').

The subset of variables selected to generate the multiple regression model are these (P.c. indicates the Person coefficient of each variable):

- Expected Years of Schooling (years) [P.c.: 0.896938]
- Mean Years of Schooling (years) [P.c.: 0.889827]
- Total Fertility Rate (live births per woman) [P.c.: -0.851443]
- Crude Birth Rate (births per 1,000 population) [P.c.: -0.874986]
- Adolescent Birth Rate (births per 1,000 women ages 15-19) [P.c.: -0.790792]
- Net Reproduction Rate (surviving daughters per woman) [P.c.: -0.761124]
- Coefficient of human inequality [P.c.: -0.703339]
- Median Age, as of 1 July (years) [P.c.: 0.786724]

The linear model obtained has the following coefficients:

Intercept: 60.230469183112646

Coefficients: [3.84139630e-01 -2.66957876e-01 -8.54299584e+00 -3.78035444e-01
-1.98550446e-02 2.50527878e+01 -6.89802117e-02 3.83722522e-01]

The resulting R-Squared is 0.8781166276189235.

The MSE is 29.899182603718497.

The new model shows a significant improvement from the previous one.

References

- [1] F. Pedregosa, G. Varoquaux, A. Gramfort, *et al.*, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [2] UNDP, *Human development index (hdi)*, UN web, Read: 2024-FEB-15, 2024. [Online]. Available: <https://hdr.undp.org/data-center/human-development-index#/indicies/HDI>.

A Source code

This is the complete listing of the source code.

A.1 Problem 1 and 2

Code for splitting the datasets into training and test sets and training a linear model.

```

1 # Reynir Siik, Franco Zambon
2 # DIT407 lp3 2024-02-11
3 # Assignment 4
4 # Problem 1, 2
5 # Life expectancy
6 # Splitting the dataset by random
7 #####
8 import numpy as np
9 import pandas as pd
10 from sklearn.model_selection import train_test_split
11 import matplotlib.pyplot as plt
12 from sklearn.linear_model import LinearRegression
13 from sklearn.metrics import r2_score
14 #####
15 filelocation_in = r"Uppgift04\00_Uppg_Data\" # folder to get data from
16 filelocation_out = r"Uppgift04\test_data\" # folder to write data to
17 ## #####
18 ## How many parts to split data into #####
19 chunks = 2
20 ## #####
21 df = pd.read_csv(filelocation_in + "life_expectancy.csv")
22 df.insert(1, 'Dataset', 0)
23 df_shape = df.shape[0]
24 df_chunk = df_shape//chunks
25 # print(df_chunk, df_shape/chunks)
26 DFtrainingSet, DFtestingSet = train_test_split(
27     df, train_size=df_chunk, random_state=42)
28 # Open a new (empty) csv file and write the data to it, skip index
29 DFtrainingSet.to_csv(
30     filelocation_out + "life_expectancy.csv", index=False)
31
32 ## Split the dataset and mark each chunk with a unique number under the
33 ## column 'Dataset'. Append data to existing file
34 for i in range(chunks-2):
35     DFtrainingSet, DFtestingSet = train_test_split(
36         DFtestingSet, train_size=df_chunk, random_state=42)
37     DFtrainingSet['Dataset'] = i + 1 # Mark the data as chunk i+1
38     # Append data chunk to csv file, no header, no index
39     DFtrainingSet.to_csv(
40         filelocation_out + "life_expectancy.csv", mode='a',
41         header=False, index=False)
42
43 DFtestingSet['Dataset'] = chunks-1
44 # Append data chunk to csv file, no header, no index
45 DFtestingSet.to_csv(filelocation_out + "life_expectancy.csv", mode='a',
46                     header=False, index=False)
47
48 ## Problem 2 #####
49 # print(DFtrainingSet.shape[0])
50 numdf = DFtrainingSet.drop(["Country", "Dataset"], axis='columns')

```

```
51 ## Calculate Pearson correlation for all columns in relation to all other
52 ## columns, two by two
53 df_corr = numdf.corr('pearson')
54 df_corr['Life_Expectancy_at_Birth,both_sexes(years)'].to_csv(filelocation_out + "co
55
56 plot_colour = '#008080'
57 plot_alpha = 0.1
58 plot_width = 8
59 plot_height = 4
60 ## Scatter plot of the life expectancy and HDI
61 plt.figure(figsize=(plot_width, plot_height))
62 # Set the range of x- and y-axis
63 plt.xlim(0.2, 1.0)
64 plt.ylim(30, 90)
65
66 plt.scatter(
67     DFtrainingSet['Human_Development_Index(value)'],
68     DFtrainingSet['Life_Expectancy_at_Birth,both_sexes(years)'],
69     color=plot_colour, alpha=plot_alpha, s=50)
70 plt.grid(True, linestyle='--', alpha=0.3)
71 plt.xlabel('Human_Development_Index\n')
72 plt.ylabel('\nLife_Expectancy')
73
74 DFtrainingSet.dropna(subset=[
75     'Human_Development_Index(value)',
76     'Life_Expectancy_at_Birth,both_sexes(years)'],
77     inplace=True)
78 DFtestingSet.dropna(subset=[
79     'Human_Development_Index(value)',
80     'Life_Expectancy_at_Birth,both_sexes(years)'],
81     inplace=True)
82 ### DF_numpy = DFtestingSet.to_numpy()
83 y_train = DFtrainingSet[
84     'Life_Expectancy_at_Birth,both_sexes(years)']
85     .to_numpy().reshape(-1, 1)
86 X_train = DFtrainingSet[
87     'Human_Development_Index(value)']
88     .to_numpy().reshape(-1, 1)
89
90 y_test = DFtestingSet[
91     'Life_Expectancy_at_Birth,both_sexes(years)']
92     .to_numpy().reshape(-1, 1)
93 X_test = DFtestingSet[
94     'Human_Development_Index(value)']
95     .to_numpy().reshape(-1, 1)
96 # print(X_test)
97 # # Create linear regression object
98 regr = LinearRegression()
99
100 # # Train the model using the training sets
101 regr.fit(X_train, y_train)
102
103 print('Coefficients:\n', regr.coef_)
104 print('Intercept:', regr.intercept_)
105 # # The mean square error
106 print("Residual sum of squares: %.2f"
107       % np.mean((regr.predict(X_train) - y_train) ** 2))
108 # Explained variance score: 1 is perfect prediction
109 print('Variance score: %.2f' % regr.score(X_train, y_train))
```

```
110
111 plt.plot(X_train, regr.predict(X_train), color='blue',
112          linewidth=3)
113 plt.savefig(filelocation_out + 'scatter_plot_problem2_train.pdf')
114 #plt.show()
115
116 plt.figure(figsize=(plot_width, plot_height))
117 # Set the range of x- and y-axis
118 plt.xlim(0.2, 1.0)
119 plt.ylim(30, 90)
120 plt.grid(True, linestyle='-', alpha=0.3)
121 plt.xlabel('Human_Development_Index\n')
122 plt.ylabel('\nPredicted_Life_Expectancy')
123
124 plt.scatter(
125     DFtestingSet['Human_Development_Index_(value)'],
126     DFtestingSet['Life_Expectancy_at_Birth,_both_sexes_(years)'],
127     color=plot_colour, alpha=plot_alpha, s=50)
128
129 plt.plot(X_test, regr.predict(X_test), color='blue',
130          linewidth=3)
131
132 plt.savefig(filelocation_out + 'scatter_plot_problem2_predict.pdf')
133
134 # Explained variance score: 1 is perfect prediction
135 print('Variance_score: %.2f' % regr.score(X_test, y_test))
136
137 # Predict the dependent variable using the fitted model
138 y_pred = regr.predict(X_test)
139
140 # Calculate the R-squared
141 r_squared = r2_score(y_test, y_pred)
142 print("R-squared:", r_squared)
143
144 # # The mean square error
145 print("Residual_sum_of_squares: %.2f"
146       % np.mean((regr.predict(X_test) - y_test) ** 2))
147
148 # The correlation between the predicted values and the target variable
149 print(pd.DataFrame(
150     {'Predict': regr.predict(X_test).T[0],
151     'Target': y_test.T[0]}).corr('pearson'))
```

A.2 Problem 3 and 4

Code for Non-linear relationship and multiple linear regression.

```
1 # Reynir Siik, Franco Zambon
2 # DIT407 lp3 2024-02-11
3 # Assignment 4
4
5 ## PROBLEM 3
6
7 import pandas as pd
8 from sklearn.model_selection import train_test_split
9 from sklearn.linear_model import LinearRegression
10 from sklearn.metrics import r2_score
11 from sklearn.metrics import mean_squared_error
```



```
12
13
14 df = pd.read_csv('C:\\Users\\ilmio\\Downloads\\life_expectancy.csv')
15 df.head(60)
16
17 df_filled = df.groupby('Country').ffill().bfill()
18 df_filled.shape
19
20 num_rows_to_select = int(len(df_filled) * 70 / 100)
21 training_df = df_filled.sample(n=num_rows_to_select, random_state=42)
22 #random_sample.shape
23 training_df.head()
24
25 # Select only numeric columns
26 numeric_df = training_df.select_dtypes(include='number')
27
28 # Compute the correlation matrix
29 correlation_matrix = numeric_df.corr()
30 correlation_matrix
31
32 human_dev_corr = correlation_matrix['Human_Development_Index_(value)']
33 most_correlated_variables = human_dev_corr.drop('Human_Development_Index_(value)').d
34 filtered_variables = most_correlated_variables[(most_correlated_variables > 0.85) | (
35 filtered_variables
36
37 filtered_variable_names = list(filtered_variables.index)
38 training_df = training_df.drop(filtered_variable_names, axis=1)
39 training_df
40
41 numeric_df = training_df.select_dtypes(include='number')
42 correlation_matrix = numeric_df.corr()
43 life_expectancy_corr = correlation_matrix['Life_Expectancy_at_Birth,_both_sexes_(years)']
44 most_correlated_variables = life_expectancy_corr.drop('Life_Expectancy_at_Birth,_both
45 most_correlated_variables
46
47 X = df_filled[['Expected_Years_of_Schooling_(years)', 'Mean_Years_of_Schooling_(years)']
48
49 y = df_filled['Life_Expectancy_at_Birth,_both_sexes_(years)']
50
51 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=
52
53 model = LinearRegression()
54
55 model.fit(X, y)
56
57 print("Intercept:", model.intercept_)
58 print("Coefficients:", model.coef_)
59
60 y_pred = model.predict(X)
61
62 r_squared = r2_score(y, y_pred)
63
64 # Print the R-squared
65 print("R-squared:", r_squared)
66
67 mse = mean_squared_error(y, y_pred)
68
69 print("Mean_Squared_Error_(MSE):", mse)
70
```

```
71
72
73
74 ## PROBLEM 4
75
76 import matplotlib.pyplot as plt
77 import seaborn as sns
78 import numpy as np
79
80 X = df_filled.drop('Life_Expectancy_at_Birth, both_sexes(years)', axis=1)
81
82 for column in X.columns:
83     # Create a scatterplot
84     sns.scatterplot(x=X[column], y=y, size=20, alpha=0.2)
85     plt.xlabel(column)
86     plt.ylabel('Life_Expectancy_at_Birth, both_sexes(years)')
87     plt.title(f'Scatterplot of {column} vs. Life_Expectancy_at_Birth')
88     plt.show()
89
90 X = df_filled['Median_Age, as of 1 July(years)']
91 X = X.to_frame()
92
93 # Check the shape of X
94 print(X.shape)
95 model.fit(X, y)
96
97 print("Intercept:", model.intercept_)
98 print("Coefficients:", model.coef_)
99
100 y_pred = model.predict(X)
101
102 r_squared = r2_score(y, y_pred)
103
104 print("R-squared:", r_squared)
105
106 log_X = np.log(X['Median_Age, as of 1 July(years)'])
107 log_X = log_X.to_frame()
108
109 model.fit(log_X, y)
110 y_pred = model.predict(log_X)
111
112 r_squared = r2_score(y, y_pred)
113
114 print("R-squared:", r_squared)
115
116 x_values = X.squeeze().values # Convert DataFrame or Series to 1D array
117 y_values = y.values
118
119 # Plot scatterplot
120 sns.scatterplot(x=x_values, y=y_values, size=20, alpha=0.2)
121 plt.xlabel('Median_Age, as of 1 July(years)')
122 plt.ylabel('Life_Expectancy_at_Birth, both_sexes(years)')
123 plt.show()
124
125 x_values = log_X.squeeze().values # Convert DataFrame or Series to 1D array
126 y_values = y.values
127
128 # Plot scatterplot
129 sns.scatterplot(x=x_values, y=y_values, size=20, alpha=0.4)
```

```
130 plt.xlabel('log(MedianAge, as of 1 July (years))')
131 plt.ylabel('Life Expectancy at Birth, both sexes (years)')
132 plt.show()
133 plt.savefig('log(MedianAge)vsLifeExpectancy.pdf', bbox_inches='tight')
134
135 X = df_filled['GrossNationalIncomePerCapita(2017PPP$)']
136 X = X.to_frame()
137
138 model.fit(X, y)
139
140 print("Intercept:", model.intercept_)
141 print("Coefficients:", model.coef_)
142
143 y_pred = model.predict(X)
144
145 r_squared = r2_score(y, y_pred)
146
147 print("R-squared:", r_squared)
148
149 log_X = np.log(X['GrossNationalIncomePerCapita(2017PPP$)'])
150 log_X = log_X.to_frame()
151 model.fit(log_X, y)
152 y_pred = model.predict(log_X)
153
154 r_squared = r2_score(y, y_pred)
155
156 print("R-squared:", r_squared)
157
158 x_values = X.squeeze().values # Convert DataFrame or Series to 1D array
159 y_values = y.values
160
161 # Plot scatterplot
162 sns.scatterplot(x=x_values, y=y_values, size=20, alpha=0.4)
163 plt.title('GNI vs. Life Expectancy Scatterplot\n')
164 plt.xlabel('\nGrossNationalIncomePerCapita(2017PPP$)')
165 plt.ylabel('Life Expectancy at Birth, both sexes (years)\n')
166
167 #plt.show()
168 plt.savefig('GNIvsLifeExpect.pdf', bbox_inches='tight')
169
170 x_values = log_X.squeeze().values # Convert DataFrame or Series to 1D array
171 y_values = y.values
172
173 # Plot scatterplot
174 sns.scatterplot(x=x_values, y=y_values, size=20, alpha=0.4)
175 plt.title('log(GNI) vs. Life Expectancy Scatterplot\n')
176 plt.xlabel('\nlog(GrossNationalIncomePerCapita(2017PPP$))')
177 plt.ylabel('Life Expectancy at Birth, both sexes (years)\n')
178 #plt.show()
179 plt.savefig('log(GNI)vsLifeExpect.pdf', bbox_inches='tight')
180
181 #####
```