

DAT565/DIT407 Assignment 4

Ola Bratt Patrick Attimont
ola.bratt@gmail.com patrickattimont@gmail.com

2024-02-xx

This paper is addressing the assignment 3 study queries within the *Introduction to Data Science & AI* course, DIT407 at the University of Gothenburg and DAT565 at Chalmers. The main source of information for this project is derived from the lectures and Skiena [1]. Assignment 4 is about correlation and linear regression.

Problem 1: Splitting the data

Problem 2: Single-variable model

The variable 'Human Development Index (value)' has the strongest pearson correlation with a coefficient of 0.92. Trained model with the following variables : Human Development Index (value) The mean squared error for is 9.45. Coefficients: [48.00427807]

Problem 3: Non-linear relationship

To explore variables with non-linear relationships, we initially assess Spearman correlation coefficients, Figure 3. Initially, we eliminate the variable exhibiting the strongest Pearson correlation, namely 'Human Development Index (value)'. Subsequently, we examine the Spearman correlation.

Our analysis reveals that 'Median Age, as of 1 July (years)' exhibits the most robust Spearman correlation. Subsequently, we develop a linear model using this variable to determine the mean squared error, which amounts to 13.58, Figure 5a.

Next, we apply transformations to the variable using logarithmic, square root, and reciprocal functions, respectively Figure 4. For each transformation, we train a linear model to ascertain the resulting mean squared error. The logarithmic transformation yields the lowest mean squared error, measuring 11.52. Figure 5b.

Before transformation, the Pearson correlation registers at 0.898, while after transformation, it increases marginally to 0.913.



Figure 1: Correlation Pearson

Problem 4: Multiple linear regression

We employ Spearman correlation analysis to identify variables with strong correlations. Utilizing a predefined threshold, we pinpoint variables with significant impact. Through experimentation with various thresholds, we determine that a threshold of 0.85 yields favorable results without excessive variable inclusion.

The identified variables include: Expected Years of Schooling, female (years), Coefficient of human inequality, Gross National Income Per Capita (2017), Median Age as of 1 July (years), Rate of Natural Change (per 1,000 population), Crude Birth Rate (births per 1,000 population), Total Fertility Rate (live births per woman), and Net Reproduction Rate (surviving daughters per woman).

Conducting a multiple linear regression analysis results in a mean squared error of 2.03. The coefficients for the model are as follows: [19, -571, 243338, 61, 1.7, -2.1, 1.9, 3].

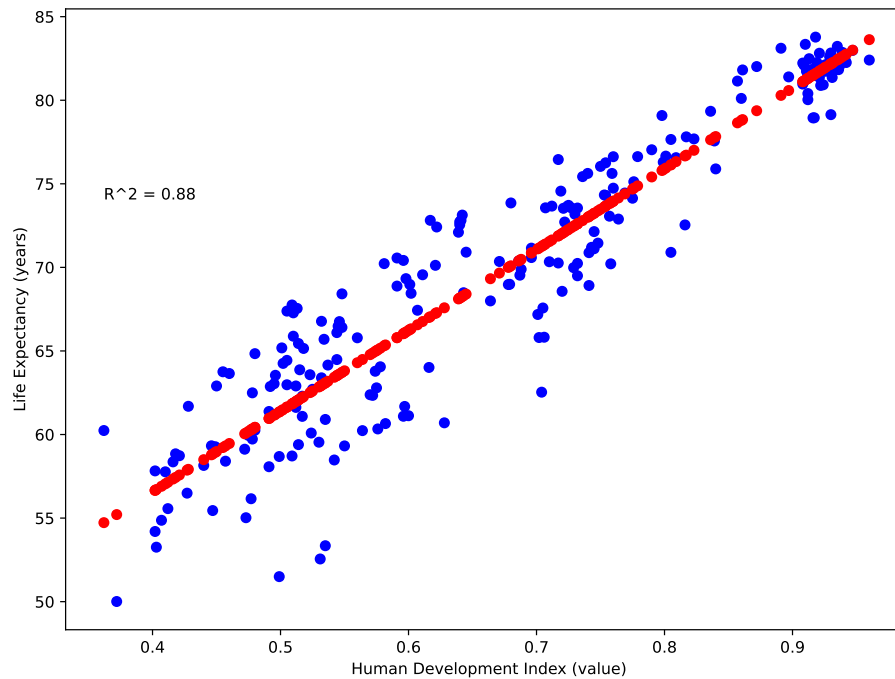


Figure 2: Linear Regression Human Development Index (value)

References

- [1] Steven S Skiena. *The Data Science Design Manual*. Retrieved 2024-01-20. 2024. URL: <https://ebookcentral.proquest.com/lib/gu/detail.action?docID=6312797>.

Appendix: Source Code

```

1 from matplotlib import pyplot
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression
6 import seaborn as sns
7 from sklearn.metrics import mean_squared_error
8 from sklearn.metrics import r2_score
9
10 def calculate_correlation(data, variable, method):
11     # Compute Pearson correlation coefficients
12     correlation_matrix = data.corr(method = method)
13
14     # Extract correlation coefficients of the target variable (life
15     #    ↳ expectancy)
16     correlation_with_life_expectancy = correlation_matrix[variable]
17     # Remove the target variable from the correlation coefficients
18     correlation_without_life_expectancy =
19     #    ↳ correlation_with_life_expectancy.drop(variable)

```

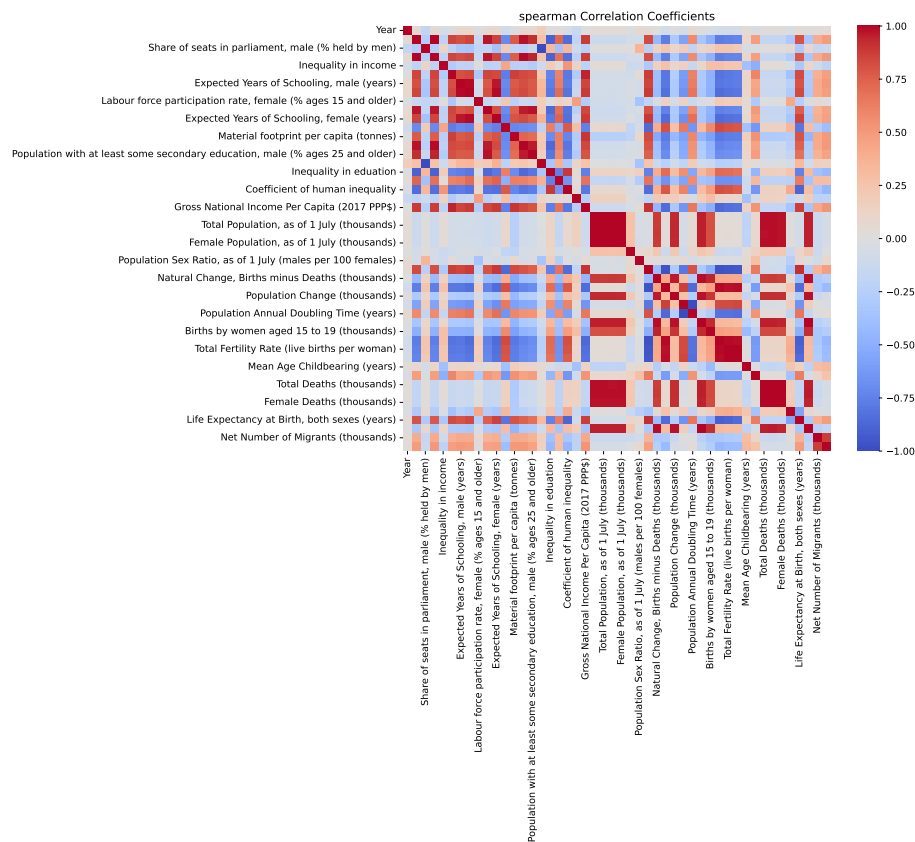


Figure 3: Correlation Spearman

```

18
19 # Find the variable with the highest absolute correlation
    ↳ coefficient
20 strongest_correlation_variable =
    ↳ correlation_without_life_expectancy.abs().idxmax()
21 strongest_correlation_coefficient =
    ↳ correlation_without_life_expectancy.abs().max()
22
23 print(f"The variable '{strongest_correlation_variable}' has the
    ↳ -strongest -" + method + f"-correlation with a-
    ↳ coefficient of {strongest_correlation_coefficient:.2f}."
    ↳ )
24
25 fig, ax = pyplot.subplots(figsize=(10, 8))
26 sns.heatmap(correlation_matrix, annot=False, cmap='coolwarm')
27 ax.set_title(method + '-Correlation-Coefficients')
28 fig.savefig(method + "_correlation.pdf", bbox_inches='tight')
29
30 return strongest_correlation_variable, correlation_matrix
31
32
33 def train_linear_regression_model(X_train, X_test, y_train, y_test,
    ↳ variables, prefix = ''):

```

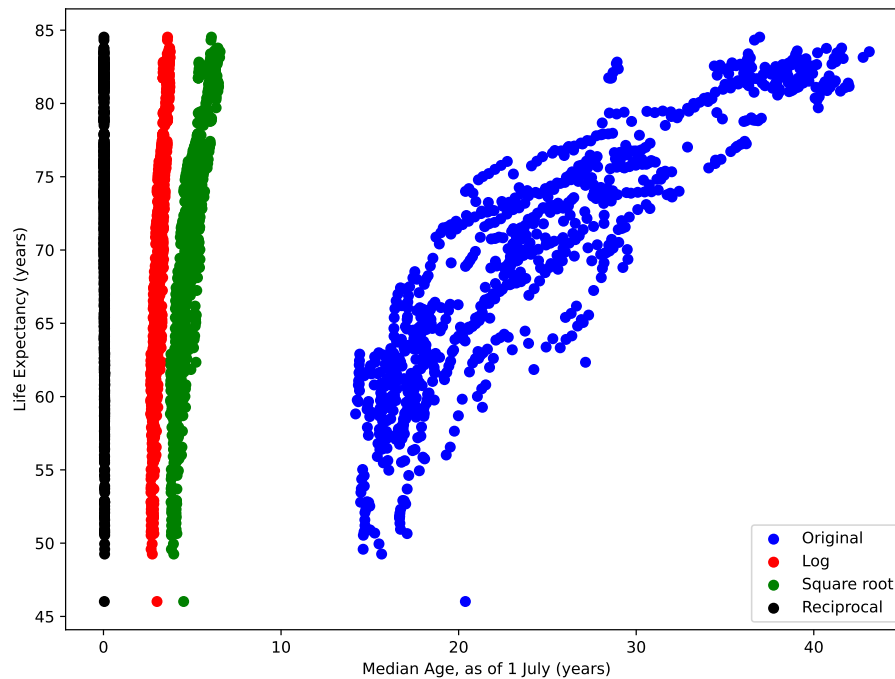
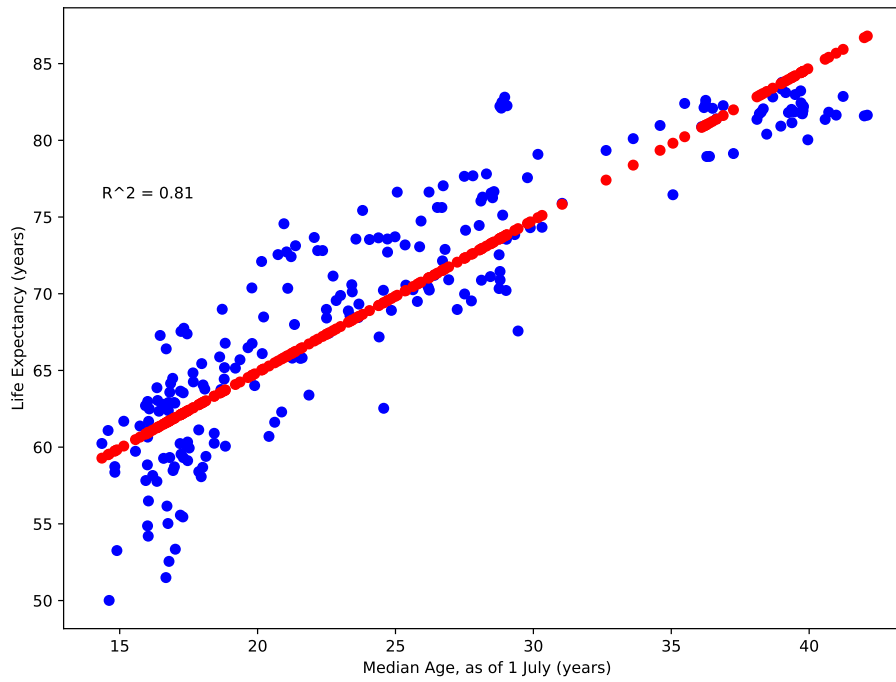


Figure 4: Linear transformation

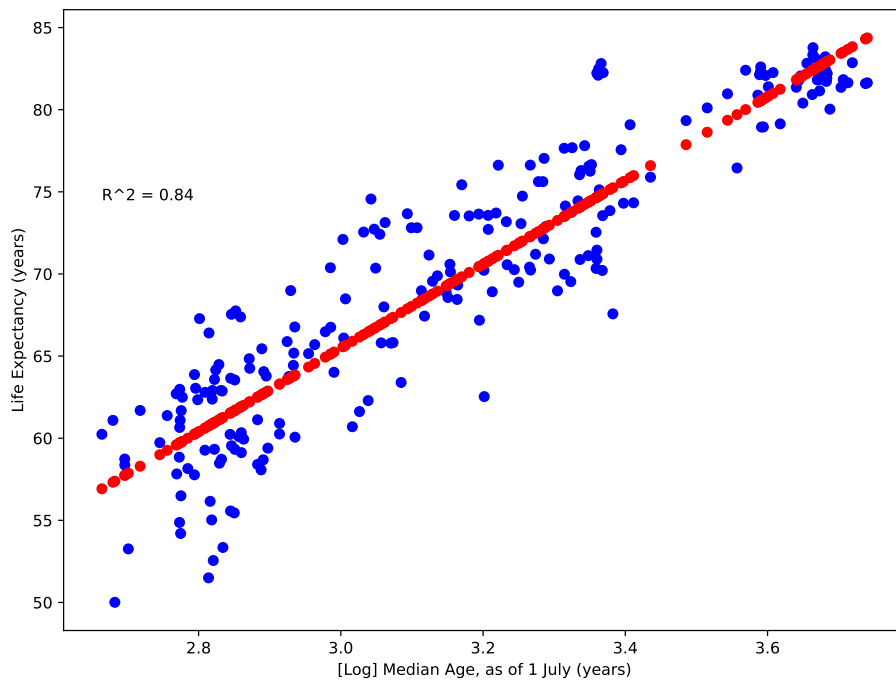
```

34     # Train a linear regression model using the variable with the
      ↳ strongest correlation
35     model = LinearRegression().fit(X_train, y_train)
36     # Make predictions
37     y_pred = model.predict(X_test)
38     r2 = r2_score(y_test, y_pred)
39     _, rows = X_test.shape
40     if rows == 1:
41         fig, ax = pyplot.subplots(figsize=(8, 6), layout='
      ↳ constrained')
42         ax.scatter(X_test, y_test, color='blue')
43         ax.scatter(X_test, y_pred, color='red')
44         ax.set_xlabel(prefix + "-" + variables)
45         ax.set_ylabel('Life Expectancy (years)')
46         ax.text(0.1, 0.7, f'R^2={r2:.2f}', ha='center', va='
      ↳ center', transform=ax.transAxes)
47         filename = prefix + "_linear_regression_" + variables + ".
      ↳ pdf"
48         filename = filename.replace('-', '_').lower()
49         fig.savefig(filename, bbox_inches='tight')
50
51
52     mse = mean_squared_error(y_test, y_pred)
53     print("Trained model with the following variables-" + prefix +
      ↳ "-:", variables)
54     print(f"The mean squared error for is-{mse:.2f}.")
55
56 def transform_variable(X_train, y_train, correlation_variable):
57     pd.options.mode.chained_assignment = None
58

```



(a) Linear Regression Median Age (original)



(b) Linear Regression Median Age (log)

Figure 5: Linear Regression Median Age

```

59     X_train_selected = X_train[[correlation_variable]]
60
61     X_train_selected['log'] = np.log(X_train[[correlation_variable
62         ↪ ]])
63     X_train_selected['sqrt'] = np.sqrt(X_train[[
64         ↪ correlation_variable]])
65     X_train_selected['reciprocal'] = 1/(X_train[[
66         ↪ correlation_variable]])
67
68     fig, ax = pyplot.subplots(figsize=(8, 6), layout='constrained')
69     ax.scatter(X_train_selected.iloc[:, 0], y_train, color='blue',
70         ↪ label='Original')
71     ax.scatter(X_train_selected['log'], y_train, color='red', label
72         ↪ ='Log')
73     ax.scatter(X_train_selected['sqrt'], y_train, color='green',
74         ↪ label='Square-root')
75     ax.scatter(X_train_selected['reciprocal'], y_train, color='
76         ↪ black', label='Reciprocal')
77     ax.set_ylabel('Life-Expectancy-(years)')
78     ax.set_xlabel(correlation_variable)
79
80     ax.legend()
81     fig.savefig("linear-transformation.pdf", bbox_inches='tight')
82
83     file_path = "../life-expectancy.csv"
84     life_expectancy = pd.read_csv(file_path, sep=',').dropna()
85     LEB = 'Life-Expectancy-at-Birth,-both-sexes-(years)'
86     life_expectancy.set_index('Country', inplace=True)
87
88     life_expectancy_train, life_expectancy_test = train_test_split(
89         ↪ life_expectancy, test_size=0.2)
90
91     X_train = life_expectancy_train.drop(LEB, axis=1)
92     X_test = life_expectancy_test.drop(LEB, axis=1)
93     y_train = life_expectancy_train[LEB]
94     y_test = life_expectancy_test[LEB]
95
96     strongest_pearson_correlation_variable, correlation_pearson =
97         ↪ calculate_correlation(life_expectancy_train, LEB, 'pearson')
98     train_linear_regression_model(X_train[[
99         ↪ strongest_pearson_correlation_variable]], X_test[[
100         ↪ strongest_pearson_correlation_variable]], y_train, y_test,
101         ↪ strongest_pearson_correlation_variable)
102
103     strongest_spearman_correlation_variable, correlation_spearman =
104         ↪ calculate_correlation(life_expectancy_train.drop(
105         ↪ strongest_pearson_correlation_variable, axis=1), LEB, '
106         ↪ spearman')
107     train_linear_regression_model(X_train[[
108         ↪ strongest_spearman_correlation_variable]], X_test[[
109         ↪ strongest_spearman_correlation_variable]], y_train, y_test,
110         ↪ strongest_spearman_correlation_variable)
111     transform_variable(X_train, y_train,
112         ↪ strongest_spearman_correlation_variable)
113     train_linear_regression_model(np.log(X_train[[
114         ↪ strongest_spearman_correlation_variable]]), np.log(X_test[[
115         ↪ strongest_spearman_correlation_variable]]), y_train, y_test,
116         ↪ strongest_spearman_correlation_variable, "[Log]")
117     train_linear_regression_model(np.sqrt(X_train[[
118         ↪ strongest_spearman_correlation_variable]]), np.sqrt(X_test[[

```

```

    ↪ strongest_spearman_correlation_variable]], y_train, y_test,
    ↪ strongest_spearman_correlation_variable, "[Sqrt]")
98 train_linear_regression_model(1/(X_train[[
    ↪ strongest_spearman_correlation_variable]], 1/(X_test[[
    ↪ strongest_spearman_correlation_variable]]), y_train, y_test,
    ↪ strongest_spearman_correlation_variable, "[Reciprocal]")
99
100 threshold = 0.85
101 correlation_spearman_no_LEB = correlation_spearman.drop([LEB])
102
103 relevant_variables = correlation_spearman_no_LEB[abs(
    ↪ correlation_spearman_no_LEB['Life-Expectancy-at-Birth', -both-
    ↪ sexes-(years)']) > threshold].index.tolist()
104 train_linear_regression_model(X_train[relevant_variables], X_test[
    ↪ relevant_variables], y_train, y_test, relevant_variables)

```