

DAT565/DIT407 Assignment 4

Ola Bratt
ola.bratt@gmail.com

Patrick Attimont
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.

Problem 3: Non-linear relationship

The variable 'Median Age, as of 1 July (years)' has the strongest spearman correlation with a coefficient of 0.91. Trained model with the following variables : Median Age, as of 1 July (years) The mean squared error for is 13.58. Trained model with the following variables [Log]: Median Age, as of 1 July (years) The mean squared error for is 11.52. Trained model with the following variables [Sqrt]: Median Age, as of 1 July (years) The mean squared error for is 12.23. Trained model with the following variables [Reciprocal]: Median Age, as of 1 July (years) The mean squared error for is 12.12.

Problem 4: Multiple linear regression

Trained model with the following variables : 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), Net Reproduction Rate (surviving daughters per woman) The mean squared error for is 2.00.

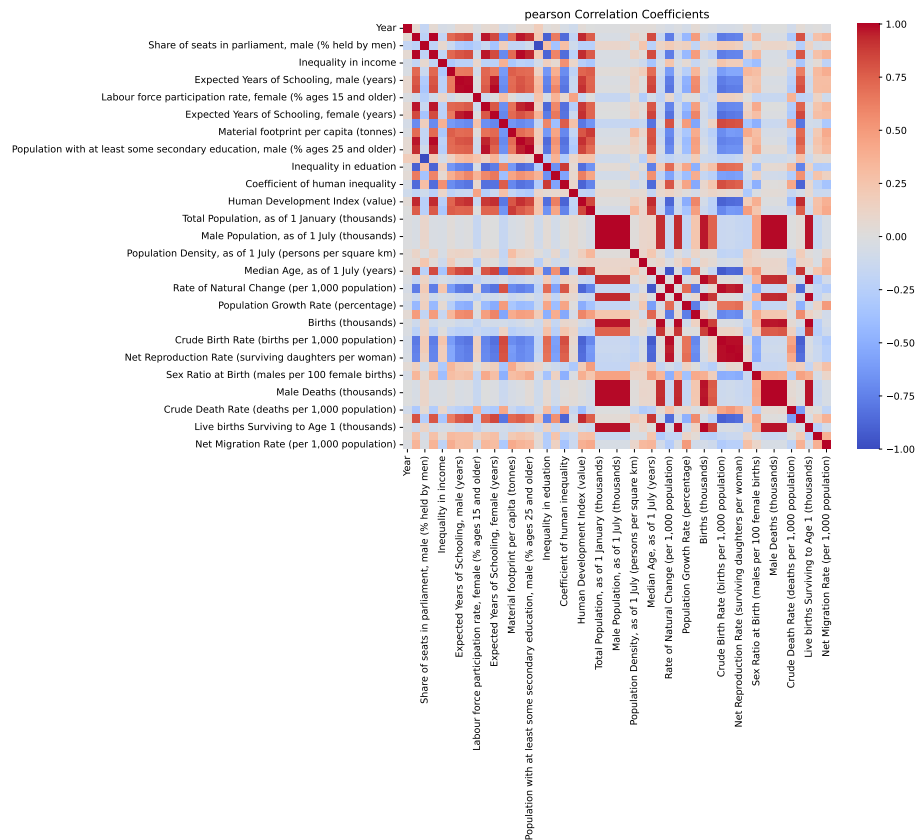


Figure 1: Correlation Pearson

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 
```

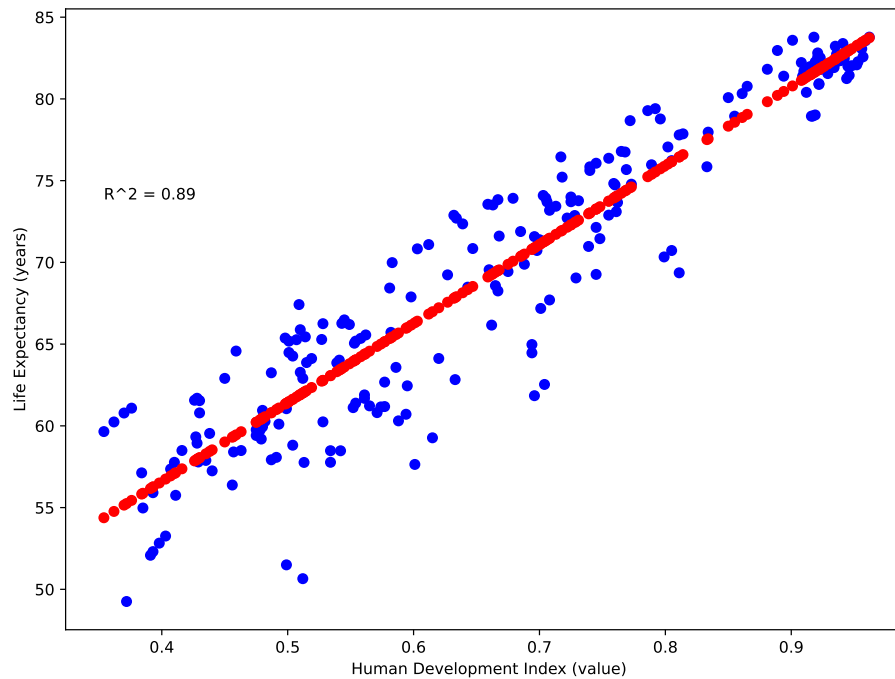


Figure 2: Linear Regression Human Development Index (value)

```

14     # Extract correlation coefficients of the target variable (life
    ↪     expectancy)
15     correlation_with_life_expectancy = correlation_matrix[variable]
16     # Remove the target variable from the correlation coefficients
17     correlation_without_life_expectancy =
    ↪     correlation_with_life_expectancy.drop(variable)
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 = ''):

```

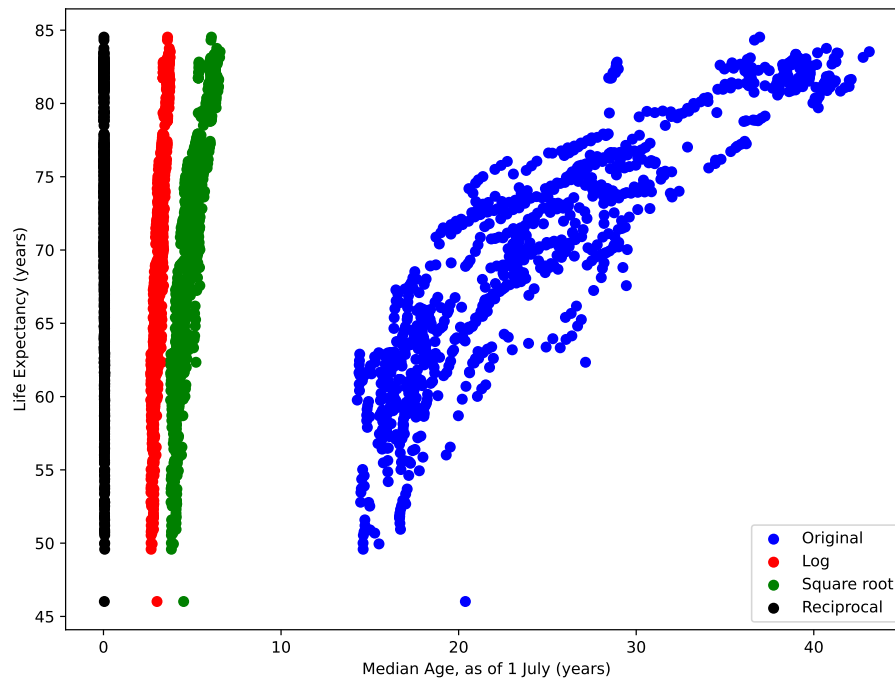
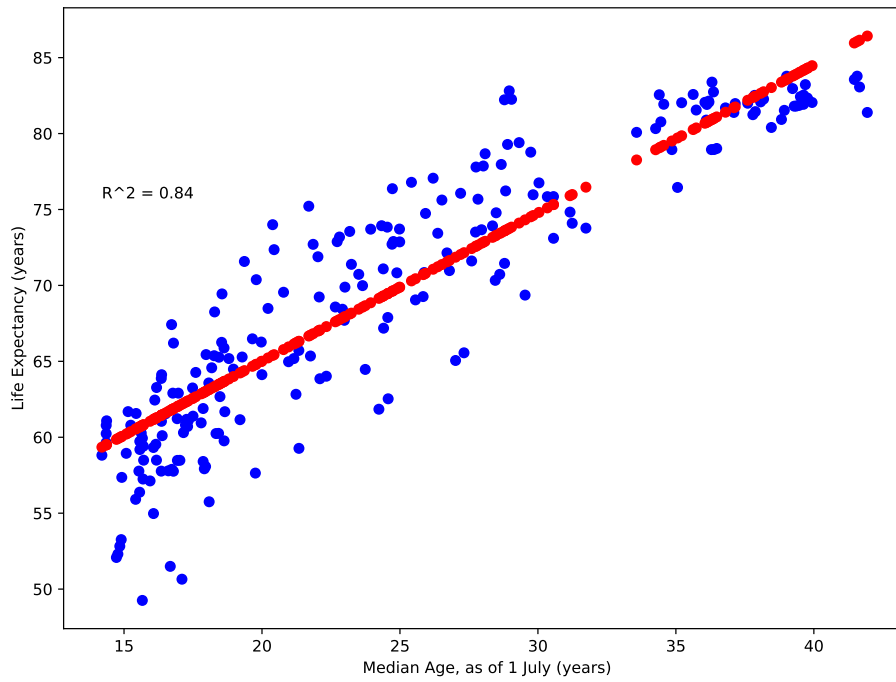


Figure 3: Linear transformation

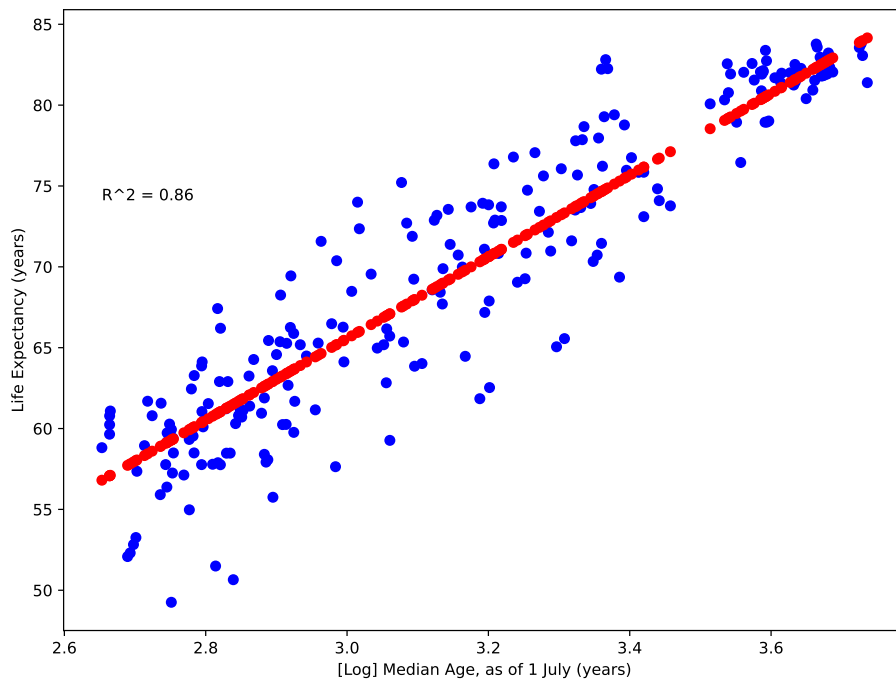
```

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 4: Linear Regression Median Age

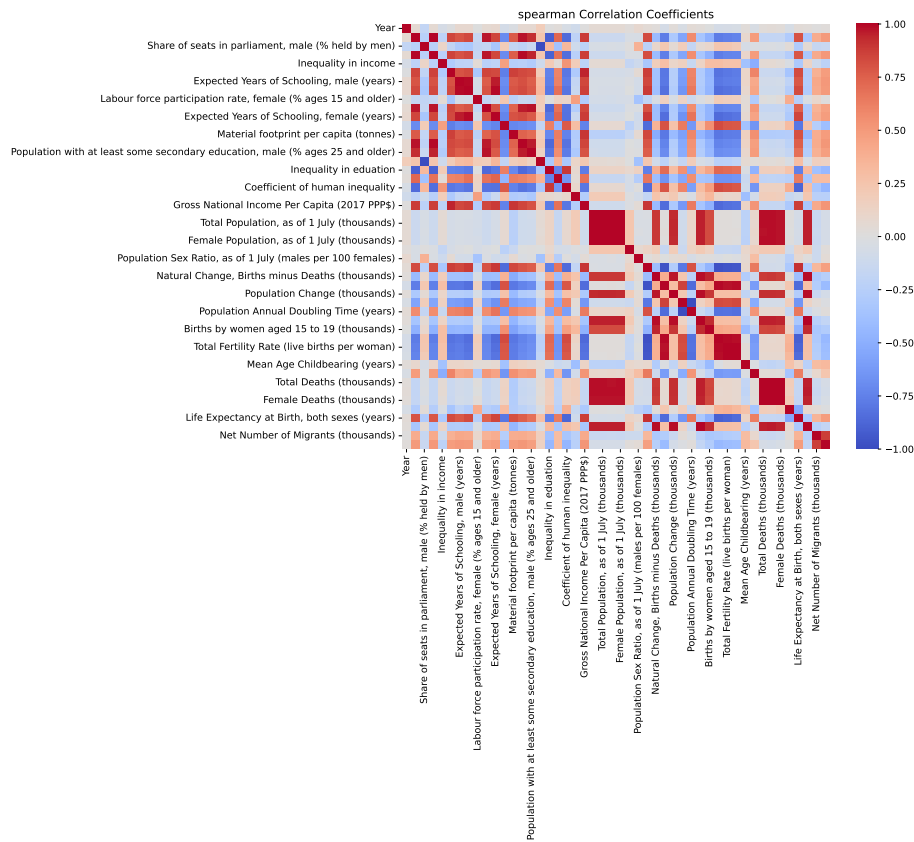


Figure 5: Correlation Spearman

```

59 X_train_selected = X_train[[ correlation_variable ]]
60
61 X_train_selected[ 'log' ] = np.log(X_train[[ correlation_variable
    ↪ ]])
62 X_train_selected[ 'sqrt' ] = np.sqrt(X_train[[
    ↪ correlation_variable ]])
63 X_train_selected[ 'reciprocal' ] = 1/(X_train[[
    ↪ correlation_variable ]])
64
65 fig, ax = pyplot.subplots(figsize=(8, 6), layout='constrained')
66 ax.scatter(X_train_selected.iloc[:, 0], y_train, color='blue',
    ↪ label='Original')
67 ax.scatter(X_train_selected[ 'log' ], y_train, color='red', label
    ↪ ='Log')
68 ax.scatter(X_train_selected[ 'sqrt' ], y_train, color='green',
    ↪ label='Square-root')
69 ax.scatter(X_train_selected[ 'reciprocal' ], y_train, color='
    ↪ black', label='Reciprocal')
70 ax.set_ylabel( 'Life-Expectancy-(years)' )
71 ax.set_xlabel( correlation_variable )
72
73 ax.legend()
74 fig.savefig("linear_transformation.pdf", bbox_inches='tight')

```

```

75
76 file_path = "../life_expectancy.csv"
77 life_expectancy = pd.read_csv(file_path, sep=',').dropna()
78 LEB = 'Life-Expectancy-at-Birth,-both-sexes-(years)'
79 life_expectancy.set_index('Country', inplace=True)
80
81 life_expectancy_train, life_expectancy_test = train_test_split(
    ↪ life_expectancy, test_size=0.2)
82
83 X_train = life_expectancy_train.drop(LEB, axis=1)
84 X_test = life_expectancy_test.drop(LEB, axis=1)
85 y_train = life_expectancy_train[LEB]
86 y_test = life_expectancy_test[LEB]
87
88
89
90 strongest_pearson_correlation_variable, correlation_pearson =
    ↪ calculate_correlation(life_expectancy_train, LEB, 'pearson')
91 train_linear_regression_model(X_train[[
    ↪ strongest_pearson_correlation_variable]], X_test[[
    ↪ strongest_pearson_correlation_variable]], y_train, y_test,
    ↪ strongest_pearson_correlation_variable)
92
93 strongest_spearman_correlation_variable, correlation_spearman =
    ↪ calculate_correlation(life_expectancy_train.drop(
    ↪ strongest_pearson_correlation_variable, axis=1), LEB, '
    ↪ spearman')
94 train_linear_regression_model(X_train[[
    ↪ strongest_spearman_correlation_variable]], X_test[[
    ↪ strongest_spearman_correlation_variable]], y_train, y_test,
    ↪ strongest_spearman_correlation_variable)
95 transform_variable(X_train, y_train,
    ↪ strongest_spearman_correlation_variable)
96 train_linear_regression_model(np.log(X_train[[
    ↪ strongest_spearman_correlation_variable]]), np.log(X_test[[
    ↪ strongest_spearman_correlation_variable]]), y_train, y_test,
    ↪ strongest_spearman_correlation_variable, "[Log]")
97 train_linear_regression_model(np.sqrt(X_train[[
    ↪ 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)

```