

Contents

1. Introduction.....	1
2. Data Preprocessing.....	1
2.1 Cleaning	1
2.2 Transformation	1
2.3 Feature Selection.....	1
2.4 Data Combination.....	2
2.5 Data Standardization.....	2
3. Develop a predictive model.....	2
3.1 Model Preparation.....	2
3.1.1 <i>Characteristic and target variable definitions</i>	2
3.1.2 <i>Divide the training set and testing set</i>	2
3.2 Model Training	2
3.3 Evaluation	3
3.4 Result.....	4
4. Conclusions	7
5. References	7

I. Introduction

The World Happiness Report has increasingly drawn global attention in recent years, as it provides valuable insights into the well - being of populations across different countries. Happiness, a complex and multi - faceted concept, is influenced by a variety of factors such as economic stability, social support, health conditions, and freedom of choice.

As we enter the mid - 2020s, understanding the trends and changes in the world's happiness levels becomes crucial. The question then arises: How have the happiness scores of countries evolved from 2024 to 2025? Are there consistent patterns among high - ranking and low - ranking countries? By analyzing the World Happiness Index maps and rankings for 2024 and 2025, we aim to explore these questions, uncover the underlying factors contributing to the differences in happiness scores, and gain a deeper understanding of the global landscape of well - being.

II. Data Preprocessing

2.1 Cleaning

There is already a happiness index and related dataset from 2018 to 2024. Considering the completeness of the data and its impact on overall statistics, the mean is used to fill in missing values.

2.2 Transformation

The data for 2022 has formatting issues (such as commas in numbers), excluding non numeric columns (such as Country and Happiness Rank), and only converting numeric columns to object types that contain numerical values.

2.3 Feature Selection

Calculate the correlation coefficient between each feature and the happiness index score. Features with higher correlation may be more important for predicting happiness index scores.

Firstly, Shapiro Wilk normality test is performed on the features of each dataset, and the Spearman correlation coefficient is selected for non normal distributions.

Then a correlation threshold (0.2) was set, retaining only features with a correlation higher than this threshold with the happiness index score, and ultimately excluding the Genealogy column.

2.4 Data Combination

Add a Year column to each dataset to indicate the year to which the data belongs. Merge all datasets into one combined.df for ease of subsequent processing.

2.5 Data Standardization

Using StandardScaler to standardize feature data, converting the standardized data into a DataFrame format, so that different features have the same scale, avoiding certain features from having excessive impact on model training due to large numerical ranges, and helping to improve model training effectiveness and convergence speed.

III. Develop a predictive model

3.1 Model Preparation

3.1.1 Characteristic and target variable definitions

Economy, Social support, Healthy life expectancy, Freedom, Perceptions of corruption, and Year are selected as characteristic variables X, which are considered to have an impact on Happiness Score. And Happiness Score itself serves as the target variable y.

3.1.2 Divide the training set and testing set

Divide the standardized feature data and target variables into training sets X_train, y_train, and testing sets X_test in an 80:20 ratio y_test Set random_state=42 during partitioning to ensure consistency of partitioning results every time the code is run, facilitating model evaluation and comparison.

3.2 Model Training

Instantiate a RandomForestRegressor model [1], set n_estimators=100 to indicate that there are 100 decision trees in the forest, and random_state=42 to ensure the repro-

ducibility of the results. Model fitting: Use the training sets X_{train} and y_{train} to fit the model, that is, let the model learn the relationship between features and target variables.

3.3 Evaluation

Predict y_{pred} on the test set X_{test} , and then calculate the mean squared error (MSE) using `mean_squared_error`, which measures the average squared error between the predicted and true values. The smaller the value, the better the model's prediction performance; Use `r2_score` to calculate the R-squared value, which represents the goodness of fit of the model to the data. The range of values is between 0 and 1, and the closer it is to 1, the better the fitting effect of the model. Print these two evaluation metrics to understand the performance of the model.

△ MSE

Mean square error is the average of the squared differences between predicted and true values, which measures the average degree of deviation between the predicted and true values. The calculation formula is as follows:

Let y_i be the true value, \hat{y}_i is the predicted value, n is the number of samples, the formula for calculating the mean square error is expressed as follows:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

△ R²

R^2 is used to measure the goodness of fit of a regression model, which represents the proportion of variance that the model can explain. The closer its value is to 1, the better the fitting effect of the model on the data.

Let \bar{y} be the average of the true values, y_i is the true value, \hat{y}_i is the predicted value. If n is the sample size, then the calculation formula for R^2 is expressed as follows:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Among them, molecule $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ represents the Sum of Squared Residuals (SSR), which is the sum of squared differences between the predicted values and the true values of the model; The denominator $\sum_{i=1}^n (y_i - \bar{y})^2$ represents the Total Sum of Squares (SST), which is the sum of squared differences between the true value and the mean of the true value.

3.4 Result

Use the trained model to predict the data for 2025 and obtain the predicted Happiness Score for each country in 2025. The results are as follows

Below are the global distribution maps of happiness index for 2024 and 2025.

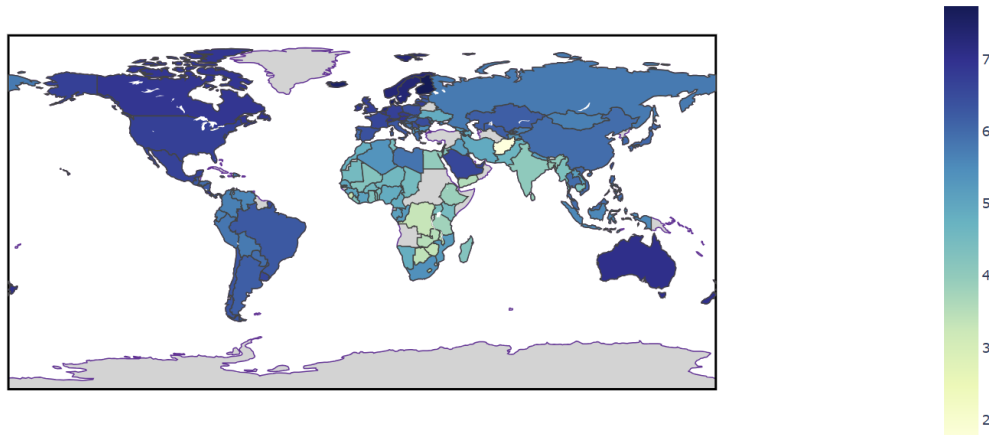


Figure 1 World Happiness Index Map in 2024

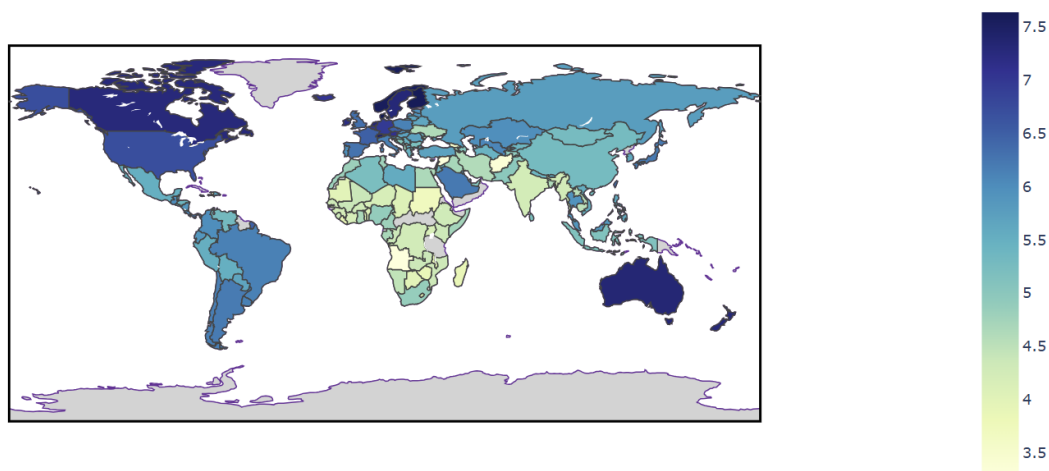


Figure 2 World Happiness Index Map in 2025

From the 2024 and 2025 World Happiness Index maps, it can be seen that the happiness index of Nordic and some Western countries such as Finland, Norway, Canada, and the United States has been at a high level for a long time, presented in dark blue to blue-green; Some African countries are mostly light green to light yellow in color, with low happiness index; There are significant differences among Asian countries. At the same time, compared with the two-year map, some countries have changed colors and their indices have fluctuated. Although the overall distribution pattern is generally

stable, local dynamic changes reflect the complex and continuous changes in factors affecting the happiness index. Countries need to constantly adjust their development strategies to ensure or enhance residents' sense of happiness.

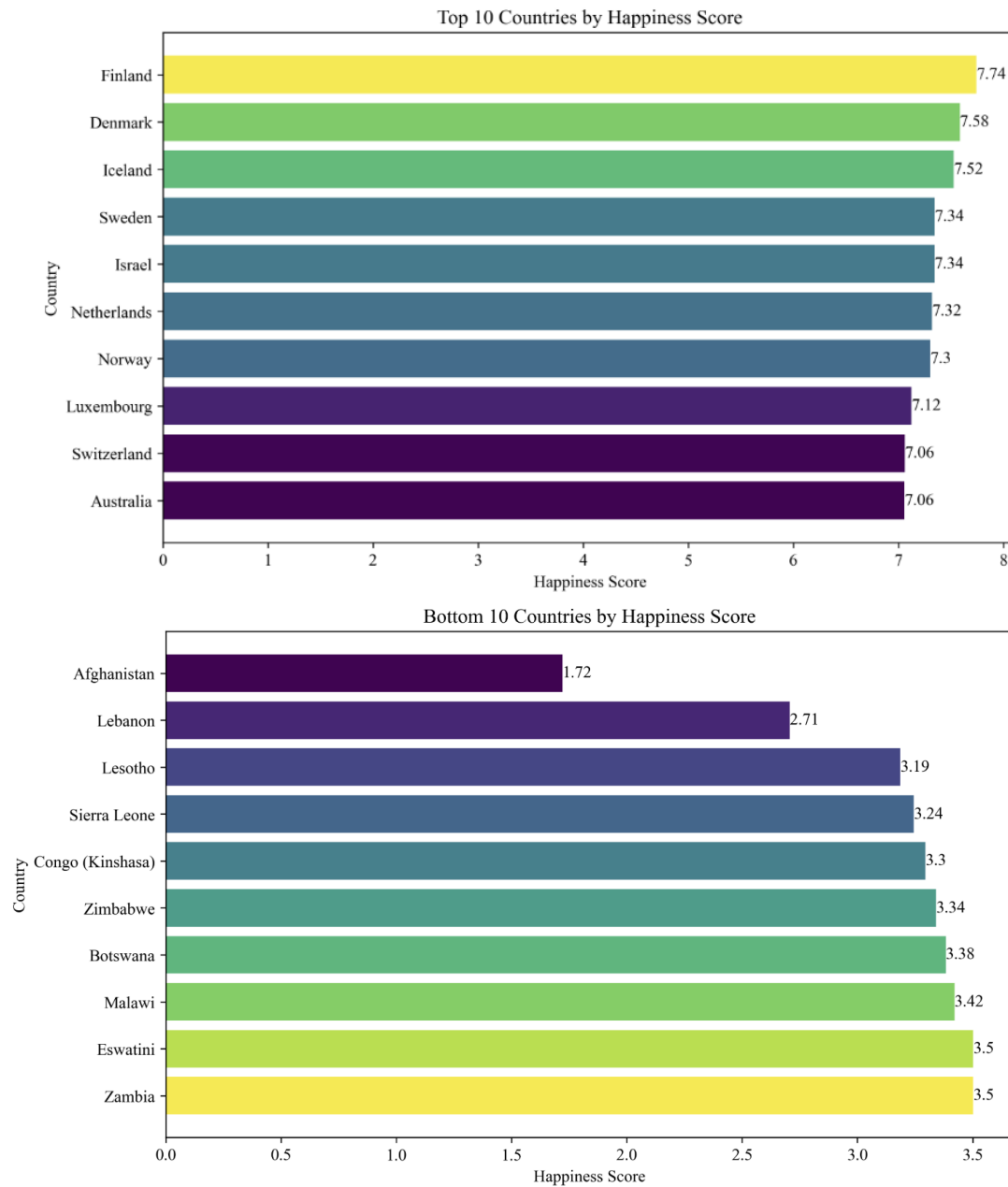


Figure 3 Happiness Rank in 2024

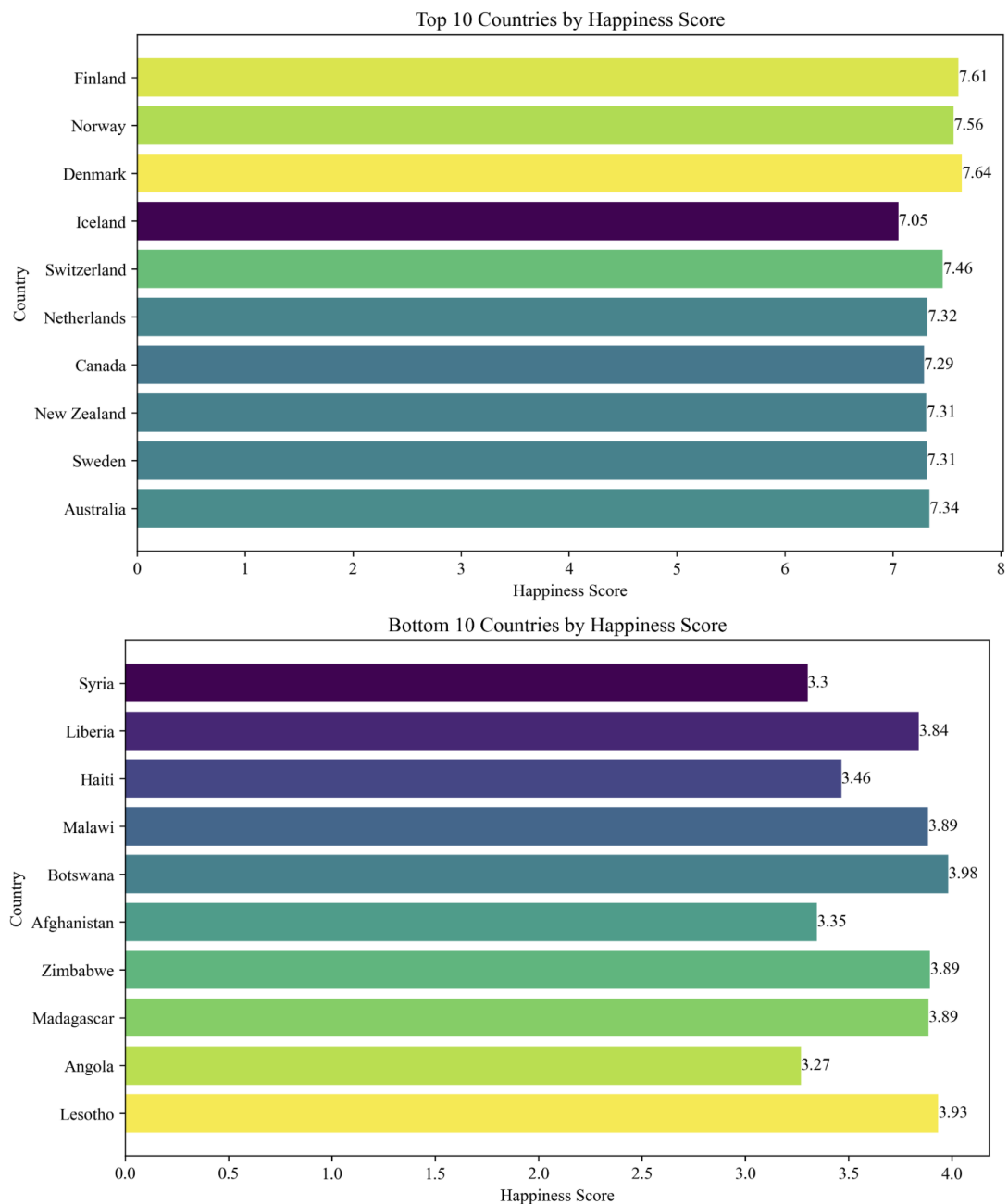


Figure 4 Happiness Rank in 2025

According to the ranking charts of happiness index in 2024 and 2025, Nordic countries such as Finland, Norway, and Denmark have consistently ranked among the top in happiness index in these two years, with high happiness scores, demonstrating their long-term advantages in improving people's happiness; Some Western countries such as Switzerland, the Netherlands, Australia, etc. are also stable in the high segment. However, countries such as Afghanistan and Syria have been in the low segment of happiness index for two years, reflecting that these regions are facing significant difficulties in social, economic, and livelihood aspects, and the task of improving people's happi-

ness is arduous. Overall, there are significant differences in happiness indices among different countries, with high happiness index countries mostly located in economically developed regions with well-established social welfare systems, while low happiness index countries are constrained by various unfavorable factors.

IV. Conclusions

According to the analysis of the predicted 2025 World Happiness Index chart and table, Nordic countries such as Finland, Norway, and Denmark rank high in happiness index, with scores mostly around 7.5 or above, while countries such as Syria and Afghanistan rank low, with scores around 3.3-3.4. The happiness score is concentrated between 3-7.6, with a larger number of countries scoring 4-6. From a regional perspective, European countries generally perform well, while Asian countries have significant differences, and African countries generally have low scores. This may be related to the situation of various countries in terms of economic development, social support, healthy life expectancy, perception of freedom and corruption. Countries with high happiness indices often perform well in these areas, while those with low happiness indices may have shortcomings.

V. References

- [1] A. Saffari, C. Leistner, J. Santner, M. Godec and H. Bischof, "On-line Random Forests," 2009 IEEE 12th International Conference on Computer Vision Workshops, ICCV Workshops, Kyoto, Japan, 2009, pp. 1393-1400, doi: 10.1109/ICCVW.2009.5457447. keywords: Machine learning;Application software;Usability;Machine learning algorithms;Computer vision;Computational efficiency;Training data;Bagging;Decision trees;Radio frequency