# Introduction to Statistic for Data Science
## Group Mini-Project Presentation: Happiness Ladder

ISDS Group 10

Christopher Barrow, Zehao Qian, Mengyuan Zhu, Hithein Augustine

Department of Natural Sciences
Durham University, England, UK

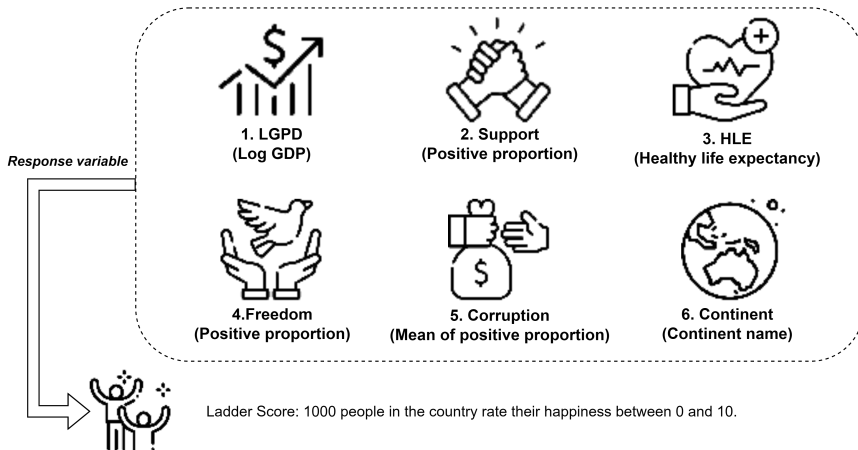December 3, 2023

# Outline

**"Happiness Ladder"**

The higher a country is on the happiness ladder, the happier, on average, its people tend to be.



Affecting factors?

**"Happiness Ladder" — Data collected from 137 countries**

*Response variable*

**1. LGPD**
**(Log GDP)**

**2. Support**
**(Positive proportion)**

**3. HLE**
**(Healthy life expectancy)**

**4.Freedom**
**(Positive proportion)**

**5. Corruption**
**(Mean of positive proportion)**

**6. Continent**
**(Continent name)**

Ladder Score: 1000 people in the country rate their happiness between 0 and 10.

# Outline

# Multiple Linear Regression Model

- **Model Target:** Some socio-economic indexes are used to predict Ladder Score to assist government decision-making.
- **Independent Variables:** LGDP, Support, HLE, Freedom, Corruption, Continent
- **Dependent Variables:** Ladder Score
- **Model Function:**
  *input*(*LGDP*, *Support*, ...) $\Rightarrow$ *output*(*Ladder Score*)
  *LadderScore* $= \beta_0 + \beta_1 * LGDP + \beta_2 * Support + ... + \epsilon$
- **Optimization Target:** Making the model with the best subset and higher Adjusted R-squared.
  $$\begin{cases} \textit{Select the independent variables} \\ \textit{Update } \beta \textit{ and } \epsilon \textit{ to minimize the residual sum of squared} \end{cases}$$

# Why we choose Multiple Linear Regression Model

| Aspect | Simple Linear Regression | Multiple Linear Regression |
|---|---|---|
| Model | $Y = \beta_0 + \beta_1 X + \varepsilon$ | $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_k X_k + \varepsilon$ |
| Advantages | Simple and easy interpretation | Captures complex relationships with multiple predictors |
| | Suitable for examining two-variable relationships | Considers multiple factors, offering a comprehensive view |
| | Less prone to overfitting with fewer predictors | Analyzes independent effects of each predictor |
| Disadvantages | Limited to two-variable relationships | More complex, challenging interpretation |
| | Assumes a linear relationship | Susceptible to multicollinearity with correlated predictors |
| | May not capture real-world complexity | More assumptions (linearity, independence, normality) |
| | | Risk of overfitting, especially with many predictors |

# Outline

# Constructing the Model – Data Processing



**Data Pre-process**

Missing Value Handling

Origin dataset → mice library → Complete dataset

Insert HLE data for State of Palestine

Country area error

Fix it with our common sense knowledge

Ecudor → South America

Uganda → Africa

**Dataset Segmentation**

Continent Mapping

"North America" = 1
"South America" = 2
"Europe" = 3
"Asia" = 4
"Africa" = 5
"Oceania" = 6

Serperate Dataset according to Continent

| Our Dataset |
|---|
| Happy_general |
| Happy_general_continent |
| Africa |
| Asia |
| Europe |
| ... |
| ... |

# Constructing the Model – Dataset Inspection

| Dataset Name | Country_name | LGDP | Support | HLE | Freedom | Corruption | Continent | Numeric Continent | Ladder_score |
|---|---|---|---|---|---|---|---|---|---|
| Happy_origin | √ | √ | √ | with NA | √ | √ | √ | | √ |
| Happy_complete | √ | √ | √ | √ | √ | √ | √ | | √ |
| Happy_general _continent | | √ | √ | √ | √ | √ | √ | √ | √ |
| Happy_general | | √ | √ | √ | √ | √ | √ | | √ |
| Africa | | √ | √ | √ | √ | √ | Africa | | √ |
| Asia | | √ | √ | √ | √ | √ | Asia | | √ |
| Europe | | √ | √ | √ | √ | √ | Europe | | √ |
| North_America | | √ | √ | √ | √ | √ | North_America | | √ |
| Oceania | | √ | √ | √ | √ | √ | Oceania | | √ |
| South_America | | √ | √ | √ | √ | √ | South_America | | √ |

# Constructing the Model – Bring dataset to MLR model



**Step 1: Correlationship**

**Step 2: Scatterplot**

**Step 3: Finding the best subset**

regsubsets()

**Step 4: Compared with**

**Single Linear Regression**

**Step 5: Regression Analytics**

Calculate the model

Data Input → Ladder_score

Model Evaluation

Variance Inflation Factor Check

Autocorrelation -- Residuals

Adjust R squared

**Future Analytics**

Interaction between variables

non-linear transformation of predictors

# Outline

# Outline