

# 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

- 1 Introduction
  - Intro to Data Set and its Context

- 2 Statistical Modeling
  - Model Explanation
  - Model Construction

- 3 Conclusions and Analysis
  - Results and Analysis
  - Conclusions

- 4 Future Work

## 1 Introduction

- Intro to Data Set and its Context

## 2 Statistical Modeling

- Model Explanation
- Model Construction

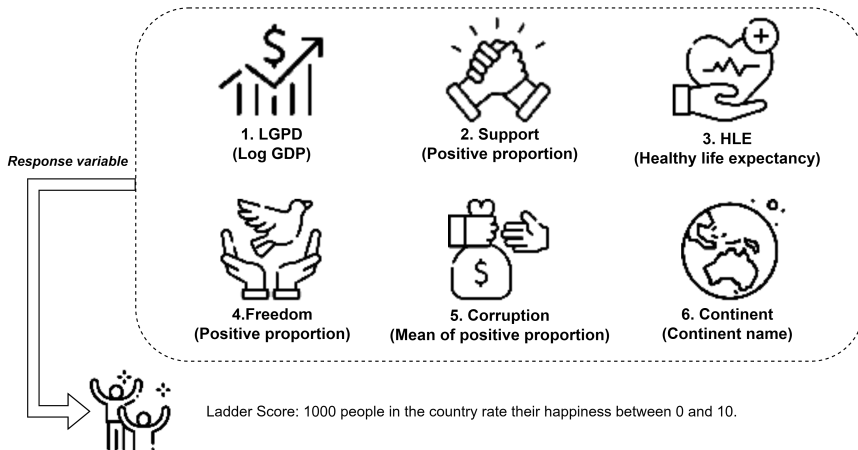
## 3 Conclusions and Analysis

- Results and Analysis
- Conclusions

## 4 Future Work

# Data Set and its Context

## “Happiness Ladder” — Data collected from 137 countries



- 1 Introduction
  - Intro to Data Set and its Context
- 2 **Statistical Modeling**
  - **Model Explanation**
  - Model Construction
- 3 Conclusions and Analysis
  - Results and Analysis
  - Conclusions
- 4 Future Work

# 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:**  
$$\text{input}(\text{LGDP}, \text{Support}, \dots) \Rightarrow \text{output}(\text{Ladder Score})$$
$$\text{LadderScore} = \beta_0 + \beta_1 * \text{LGDP} + \beta_2 * \text{Support} + \dots + \epsilon$$
- **Optimization Target:** Making the model with the **best subset** and **higher Adjusted R-squared**.  
$$\left\{ \begin{array}{l} \text{Select the independent variables} \\ \text{Update } \beta \text{ and } \epsilon \text{ to minimize the residual sum of squared} \end{array} \right.$$

# Why we choose Multiple Linear Regression Model

Aspect	Simple Linear Regression	Multiple Linear Regression
Model	$Y = \beta_0 + \beta_1 X + \varepsilon$	$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon$
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

## 1 Introduction

- Intro to Data Set and its Context

## 2 Statistical Modeling

- Model Explanation
- **Model Construction**

## 3 Conclusions and Analysis

- Results and Analysis
- Conclusions

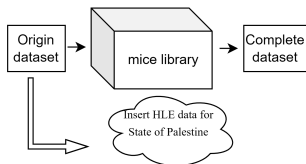
## 4 Future Work



# Constructing the Model – Data Processing

## Data Pre-process

### Missing Value Handling



### 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		√

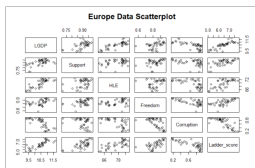
# Constructing the Model – Bring dataset to MLR model

## Correlationship

	GDP	Support	HLE	Freedom	Corruption	Ladder score
GDP	1.000000	0.738068	0.811583	0.451439	-0.430908	0.784367
Support	0.738068	1.000000	0.706043	0.541631	-0.772494	0.834517
HLE	0.811583	0.706043	1.000000	0.407733	-0.399679	0.741648
Freedom	0.451439	0.541631	0.407733	1.000000	-0.307863	0.667044
Corruption	-0.430908	-0.772494	-0.399679	-0.307863	1.000000	-0.471935
Ladder score	0.784367	0.834517	0.741648	0.667044	-0.471935	1.000000

	GDP	Support	HLE	Freedom	Corruption	Ladder score
GDP	1.000000	0.540802	0.809194	0.422627	-0.747429	0.822464
Support	0.540802	1.000000	0.361543	0.411018	-0.339348	0.659029
HLE	0.809194	0.361543	1.000000	0.363757	-0.590584	0.673260
Freedom	0.422627	0.411018	0.363757	1.000000	-0.411374	0.615077
Corruption	-0.747429	-0.339348	-0.590584	-0.411374	1.000000	-0.761485
Ladder score	0.822464	0.659029	0.673260	0.615077	-0.761485	1.000000

## Scatterplot



## Finding the best subset

## Regression Analytics

## 1 Introduction

- Intro to Data Set and its Context

## 2 Statistical Modeling

- Model Explanation
- Model Construction

## 3 Conclusions and Analysis

- **Results and Analysis**
- Conclusions

## 4 Future Work

# Model Results and Analysis

## 1 Introduction

- Intro to Data Set and its Context

## 2 Statistical Modeling

- Model Explanation
- Model Construction

## 3 Conclusions and Analysis

- Results and Analysis
- **Conclusions**

## 4 Future Work

# Conclusions and Insights

# Next Steps