

# Group 12 MAT3024 REGRESSION ANALYSIS Assignment

Group 12

2024-07-15

## 1. Topic: Identifying Key Factors Influencing the Average Score in Global Well-being Indices

### Introduction

Measuring and comprehending the state of the world's well-being has been a top priority for international organisations, researchers, and governments in recent decades. The notion of well-being extends beyond economic metrics and encompasses a comprehensive assessment of the physical, mental, and social aspects of people's and society's quality of life. Global well-being indices, including the Human Development Index and the World Happiness Report, are examples of instruments that play a significant role in evaluating and comparing the levels of well-being among countries.

### 2. Problem Identification and Objectives

These indices are derived from a combination of various indicators, including economic, social, environmental, and governance factors. However, understanding which specific factors most significantly influence the average score of this index can be challenging due to the complexity and interrelationships among the indicators.

The goal of this study is to examine the wide range of variables that affect the average score on global well-being indicators. Policymakers can better develop policies to improve the general well-being and quality of life for populations worldwide by identifying and analysing these elements. By doing so, policymakers and stakeholders can prioritize resources and interventions to improve overall well-being more effectively. The specific objectives are to:

1. **Determine the correlation between different indicators and the average score.**
2. **Develop multiple regression models to identify significant predictors of the average score.**
3. **Compare the models using various statistical criteria (AIC, BIC, Adjusted  $R^2$ , Mallow's  $C_p$ ) to find the best fitting model.**
4. **Provide actionable insights based on the final model to guide policy and decision-making.**

## 3. Data

### Data Information and Source

**Data Source:** The dataset used in this analysis is sourced from Kaggle, specifically from the dataset titled "2023 Global Country Development and Prosperity Index," which is available at the following URL: [2023 Global Country Development and Prosperity Index](#).

**Dataset Description:** This dataset provides comprehensive information on various indicators that measure the development and prosperity of countries globally. It includes data for the year 2023 and covers a wide range of factors that contribute to a country's overall well-being and development.

**Key Variables:**

1. **AveragScore:** The overall score representing a country's development and prosperity.
2. **InvestmentEnvironment:** Measures the conduciveness of the environment for investment activities.
3. **Education:** Represents the quality and accessibility of the education system.
4. **PersonelFreedom:** Indicates the level of personal freedom experienced by the citizens.
5. **SafetySecurity:** Reflects the safety and security conditions within the country.
6. **SocialCapital:** Represents the strength of social networks and community engagement.
7. **MarketAccessInfrastructure:** Measures the accessibility and quality of market infrastructure.
8. **Governance:** Reflects the effectiveness and quality of governance structures.
9. **EconomicQuality:** Represents the overall economic stability and quality.
10. **Health:** Measures the quality and accessibility of healthcare services.
11. **NaturalEnvironment:** Represents the quality of the natural environment.
12. **LivingConditions:** Reflects the general living conditions within the country.
13. **EnterpriseConditions:** Measures the conduciveness of the environment for enterprise activities.

**Data Collection Method:** The data in this dataset is likely collected from various reputable sources and international organizations that track development indicators. These sources may include government reports, international agencies, non-governmental organizations, and research institutions.

**Purpose of the Dataset:** The primary purpose of this dataset is to provide a comprehensive view of the factors contributing to the development and prosperity of countries. It can be used for comparative analysis, policy-making, academic research, and understanding the areas that need improvement for enhancing the overall well-being of populations globally.

## 4. Analysis

**Explain the relationship between X variable and Y variable**

### Explanation of the Choice of X and Y Variables

In our analysis, the y-variable (dependent variable) is **AveragScore**, which represents the average well-being score of a country. The choice of this variable is driven by the objective to understand and quantify the factors that contribute to the overall well-being of nations. The well-being score encapsulates various aspects of life quality, making it a comprehensive measure of societal health and prosperity.

The x-variables (predictors) selected for the analysis are:

#### 1. Investment Environment:

- **Reason for Selection:** The investment environment is critical as it reflects the economic opportunities and stability within a country. A favorable investment environment often leads to economic growth and improved living standards.

#### 2. Education:

- **Reason for Selection:** Education is a fundamental driver of individual and societal development. Higher education levels are associated with better employment opportunities, higher incomes, and improved health outcomes, all of which contribute to overall well-being.

### 3. Personal Freedom:

- **Reason for Selection:** Personal freedom encompasses civil liberties and political rights. Societies with higher levels of personal freedom tend to have happier citizens due to greater autonomy, better self-expression, and participation in civic activities.

### 4. Safety and Security:

- **Reason for Selection:** Safety and security are essential components of well-being. High levels of crime and violence can significantly diminish quality of life and overall happiness.

### 5. Social Capital:

- **Reason for Selection:** Social capital refers to the networks and relationships that facilitate collective action and community support. Strong social capital can lead to better health outcomes, reduced crime rates, and increased civic engagement.

### 6. Market Access and Infrastructure:

- **Reason for Selection:** Market access and infrastructure are vital for economic activities, including trade and commerce. Good infrastructure supports economic growth and accessibility, enhancing the quality of life.

### 7. Governance:

- **Reason for Selection:** Good governance, including transparency, accountability, and effective institutions, is crucial for ensuring equitable resource distribution and maintaining public trust.

### 8. Economic Quality:

- **Reason for Selection:** Economic quality includes factors like income distribution, economic stability, and productivity. These elements directly affect individuals' living conditions and opportunities.

### 9. Health:

- **Reason for Selection:** Health is a fundamental aspect of well-being. Access to healthcare services and overall population health status are crucial for maintaining a high quality of life.

### 10. Natural Environment:

- **Reason for Selection:** The quality of the natural environment, including air and water quality, biodiversity, and green spaces, significantly impacts physical and mental health.

### 11. Living Conditions:

- **Reason for Selection:** Living conditions encompass housing quality, access to basic services, and overall comfort. Good living conditions are directly linked to higher life satisfaction.

### 12. Enterprise Conditions:

- **Reason for Selection:** Enterprise conditions include the ease of doing business, innovation, and entrepreneurship. A thriving business environment fosters job creation and economic growth.

## 2. Justification of X Variable Selection

The selection of the predictor variables (X variables) for analyzing their relationship with the average well-being score (Y variable) is based on extensive research and understanding of the factors that significantly impact societal well-being. Here's the justification for including each predictor:

## 1. Investment Environment

- **Rationale:** A favorable investment environment attracts foreign direct investment (FDI) and domestic investments, which can lead to economic growth, job creation, and infrastructure development (Dunning, 2002).
- **Impact:** Enhanced economic opportunities and improved infrastructure can elevate the average well-being score of a nation by increasing employment, income levels, and access to essential services (Dunning, 2002).

## 2. Social Capital

- **Rationale:** Social capital refers to the networks, norms, and trust that facilitate coordination and cooperation among people (Putnam, 2000).
- **Impact:** Strong social networks and community bonds provide emotional support, improve mental health, and foster a sense of belonging, contributing positively to well-being (Putnam, 2000).

## 3. Market Access and Infrastructure

- **Rationale:** Effective infrastructure and market access are crucial for economic activities, connectivity, and service delivery (Calderón & Servén, 2004).
- **Impact:** Good infrastructure and easy market access improve living conditions by facilitating trade, transportation, and access to services, thereby enhancing well-being (Calderón & Servén, 2004).

## 4. Education

- **Rationale:** Education is a cornerstone of personal and professional development. It equips individuals with knowledge, skills, and opportunities, leading to improved life outcomes (Hanushek & Woessmann, 2010).
- **Impact:** Higher education levels are strongly associated with better employment prospects, higher income, and enhanced social mobility, contributing to individual and societal well-being (Hanushek & Woessmann, 2010).

## 5. Personal Freedom

- **Rationale:** Personal freedom, encompassing civil liberties and political rights, is fundamental to human dignity and autonomy (Sen, 1999).
- **Impact:** Societies that protect and promote personal freedoms tend to have higher levels of happiness and life satisfaction, as individuals can freely pursue their goals and aspirations (Sen, 1999).

## 6. Safety and Security

- **Rationale:** Physical safety and security are essential for a stable and peaceful society. Without these, individuals' well-being is significantly compromised (Wilkinson & Pickett, 2009).
- **Impact:** High levels of safety and security reduce fear and stress, allowing individuals to lead more productive and satisfying lives (Wilkinson & Pickett, 2009).

## 7. Governance

- **Rationale:** Good governance ensures the fair and efficient management of resources, transparency, and accountability in public affairs (Kaufmann, Kraay, & Mastruzzi, 2009).
- **Impact:** Countries with effective governance structures tend to have higher public trust, better public services, and reduced corruption, all of which contribute to societal well-being (Kaufmann et al., 2009).

## 8. Economic Quality

- **Rationale:** Economic quality reflects the overall health of the economy, including factors like income distribution, economic stability, and employment (Stiglitz, Sen, & Fitoussi, 2009).
- **Impact:** A high-quality economy supports sustainable development, equitable wealth distribution, and economic resilience, positively impacting citizens' well-being (Stiglitz et al., 2009).

## 9. Health

- **Rationale:** Health is a fundamental aspect of human life and well-being. Access to healthcare services and overall health status are critical determinants of quality of life (Marmot & Wilkinson, 2005).
- **Impact:** Good health enables individuals to lead productive lives, reduces healthcare costs, and improves life expectancy, directly influencing well-being (Marmot & Wilkinson, 2005).

## 10. Natural Environment

- **Rationale:** The quality of the natural environment, including air and water quality, biodiversity, and green spaces, significantly affects physical and mental health (McMichael, Woodruff, & Hales, 2006).
- **Impact:** A healthy natural environment supports physical health, reduces stress, and provides recreational opportunities, enhancing overall well-being (McMichael et al., 2006).

## 11. Living Conditions

- **Rationale:** Adequate living conditions, including housing quality, access to clean water, sanitation, and electricity, are basic human needs (United Nations, 2015).
- **Impact:** Improved living conditions lead to better health outcomes, reduced poverty, and higher life satisfaction, contributing to well-being (United Nations, 2015).

## 12. Enterprise Conditions

- **Rationale:** A supportive business environment encourages entrepreneurship, innovation, and economic diversification (Naudé, 2010).
- **Impact:** Favorable enterprise conditions lead to job creation, higher incomes, and economic growth, which are crucial for improving well-being (Naudé, 2010).

## Summary

The chosen X variables represent a comprehensive set of factors that collectively influence the average well-being score (Y variable). Each predictor has a theoretical and empirical basis for its inclusion, ensuring a holistic analysis of the determinants of well-being. The subsequent regression analysis confirms the significance and relative importance of these predictors in explaining variations in well-being scores across different countries.

```
# Load necessary libraries
library(readxl)
library(ggplot2)
library(leaps)
library(car)
```

```
## Loading required package: carData
```

```
library(MASS)
library(corrplot)
```

```
## corrplot 0.92 loaded
```

```
library(olsrr)
```

```
##
## Attaching package: 'olsrr'

## The following object is masked from 'package:MASS':
##
##      cement

## The following object is masked from 'package:datasets':
##
##      rivers
```

```
library(AICcmodavg)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following object is masked from 'package:MASS':
##
##      select

## The following object is masked from 'package:car':
##
##      recode

## The following objects are masked from 'package:stats':
##
##      filter, lag

## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union
```

```
library(knitr)
library(kableExtra)
```

```
##
## Attaching package: 'kableExtra'

## The following object is masked from 'package:dplyr':
##
##      group_rows
```

```

library(caret)

## Loading required package: lattice

library(reshape2)
library(gridExtra)

##
## Attaching package: 'gridExtra'

## The following object is masked from 'package:dplyr':
##
##      combine

# Load the dataset
file_path <- "C:/Users/wanda/Downloads/archive (8)/data.csv" # Update this to your actual file path
data_new <- read.csv(file_path)
# Inspect the structure of the dataset
str(data_new)

## 'data.frame':    167 obs. of  14 variables:
##  $ Country      : chr  " Denmark" " Sweden" " Norway" " Finland" ...
##  $ AveragScore   : num  84.5 83.7 83.6 83.5 83.4 ...
##  $ SafetySecurity : num  92.6 91 93.3 89.6 95.7 ...
##  $ PersonelFreedom : num  94.1 91.9 94.1 92 87.5 ...
##  $ Governance     : num  89.5 86.4 89.7 90.4 87.7 ...
##  $ SocialCapital  : num  82.6 78.3 79 77.3 69.1 ...
##  $ InvestmentEnvironment : num  82.4 82.8 82.2 84.1 80.8 ...
##  $ EnterpriseConditions : num  79.6 75.5 76 77.2 83.8 ...
##  $ MarketAccessInfrastructure: num  78.8 79.7 75.9 78.8 78.7 ...
##  $ EconomicQuality : num  76.8 76.2 77.2 70.3 79.7 ...
##  $ LivingConditions : num  95.8 95.3 94.7 94.5 94.7 ...
##  $ Health         : num  81.1 82.3 83 81.2 82.1 ...
##  $ Education      : num  87.5 85.9 85.7 88.4 87.7 ...
##  $ NaturalEnvironment : num  73.9 78.7 72.4 78 73.6 ...

# Remove rows with missing 'AveragScore' values
data_new_cleaned <- na.omit(data_new)

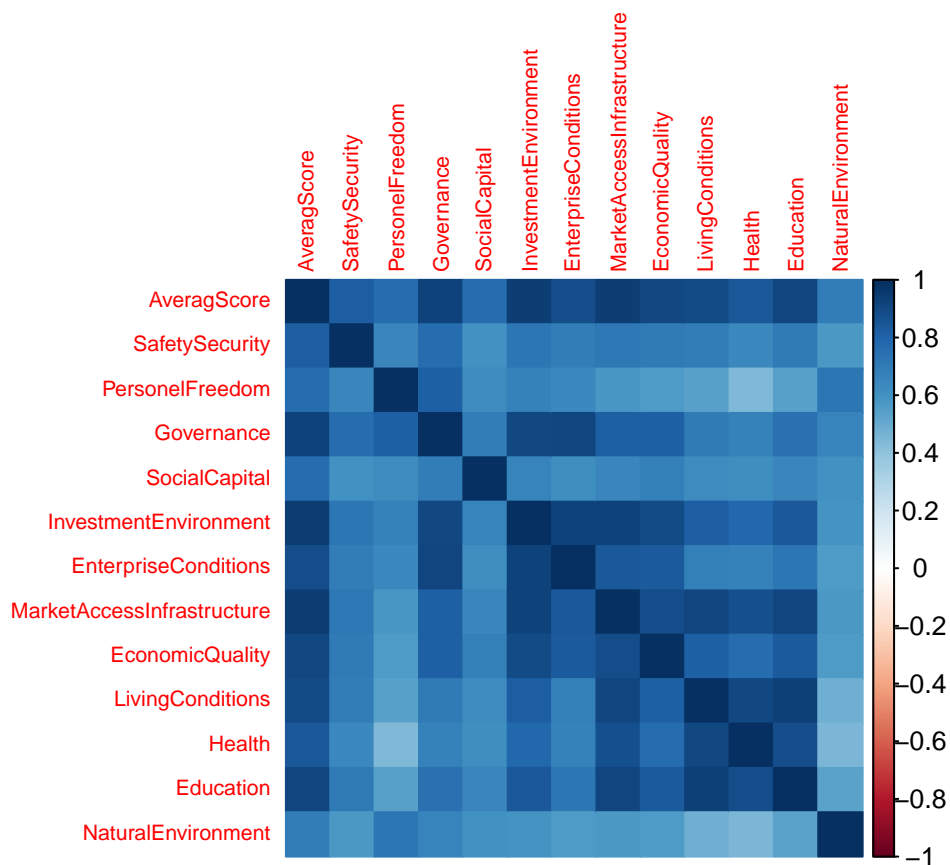
# Correlation analysis
numerical_columns <- names(data_new_cleaned)[sapply(data_new_cleaned, is.numeric)]
data_selected <- data_new_cleaned[, numerical_columns]
correlation_matrix <- cor(data_selected)

# Print the correlation matrix
correlation_matrix %>%
  round(2) %>%
  kable(caption = "Correlation Matrix") %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed", "responsive"))

```

	AveragScore	SafetySecurity	PersonelFreedom	Governance	SocialCapital	Investmen
AveragScore	1.00	0.83	0.76	0.92	0.76	
SafetySecurity	0.83	1.00	0.66	0.76	0.61	
PersonelFreedom	0.76	0.66	1.00	0.82	0.63	
Governance	0.92	0.76	0.82	1.00	0.70	
SocialCapital	0.76	0.61	0.63	0.70	1.00	
InvestmentEnvironment	0.95	0.73	0.67	0.91	0.66	
EnterpriseConditions	0.88	0.69	0.65	0.91	0.61	
MarketAccessInfrastructure	0.94	0.71	0.59	0.81	0.66	
EconomicQuality	0.91	0.71	0.57	0.82	0.68	
LivingConditions	0.89	0.70	0.54	0.71	0.63	
Health	0.85	0.64	0.45	0.67	0.62	
Education	0.91	0.71	0.55	0.74	0.66	
NaturalEnvironment	0.69	0.58	0.72	0.66	0.60	

```
# Visualize the correlation matrix
corrplot(correlation_matrix, method = "color", tl.cex = 0.7)
```



The correlation matrix shown in the image provides the Pearson correlation coefficients between each pair of variables in the dataset. The values range from -1 to 1, where:

- 1 indicates a perfect positive correlation.



- -1 indicates a perfect negative correlation.
- 0 indicates no correlation.

Here's an interpretation of the key relationships between variables:

### Key Correlations with AverageScore

- **SafetySecurity (0.83)**: There is a strong positive correlation between AverageScore and SafetySecurity. This suggests that countries with higher safety and security tend to have higher average scores.
- **Governance (0.92)**: This is one of the highest correlations, indicating that better governance is strongly associated with higher average scores.
- **InvestmentEnvironment (0.95)**: This is the highest correlation with AverageScore, suggesting that a favorable investment environment is crucial for higher average scores.
- **MarketAccessInfrastructure (0.94)**: Another strong positive correlation, indicating that better market access and infrastructure are associated with higher average scores.
- **EconomicQuality (0.91)**: Strongly positive, suggesting that higher economic quality is associated with higher average scores.
- **LivingConditions (0.89)**: Indicates that better living conditions are associated with higher average scores.
- **Health (0.85)**: Shows a strong positive correlation, indicating that better health conditions are associated with higher average scores.
- **Education (0.91)**: Another strong positive correlation, suggesting that better education systems are associated with higher average scores.

### Other Notable Correlations

- **SafetySecurity and Governance (0.76)**: Good governance is often associated with better safety and security in a country.
- **SafetySecurity and InvestmentEnvironment (0.73)**: Safer countries tend to have better investment environments.
- **Governance and InvestmentEnvironment (0.91)**: Indicates that countries with good governance also tend to have favorable investment environments.
- **MarketAccessInfrastructure and InvestmentEnvironment (0.93)**: Better market access and infrastructure are associated with a better investment environment.
- **EconomicQuality and MarketAccessInfrastructure (0.92)**: High economic quality is associated with good market access and infrastructure.
- **Education and LivingConditions (0.94)**: Better education systems are strongly correlated with better living conditions.
- **Health and Education (0.88)**: Good health conditions are strongly correlated with better education systems.

## Low or Negative Correlations

- **SocialCapital and Health (0.67):** The correlation is relatively lower compared to others, suggesting that social capital is less associated with health conditions.
- **SocialCapital and InvestmentEnvironment (0.66):** Social capital is less strongly correlated with the investment environment compared to other variables.

## Summary

The correlation matrix indicates that several factors are strongly associated with AverageScore. In particular, InvestmentEnvironment, Governance, MarketAccessInfrastructure, EconomicQuality, and Education show very high positive correlations with AverageScore. This suggests that improvements in these areas are likely to be associated with higher overall scores.

These correlations help in understanding the relationships between different indicators and can guide policymakers in focusing on the most impactful areas to improve overall scores.

```
# Simple linear regression models
simple_model1 <- lm(AveragScore ~ LivingConditions, data = data_selected)
simple_model2 <- lm(AveragScore ~ Health, data = data_selected)
simple_model3 <- lm(AveragScore ~ Education, data = data_selected)
simple_model4 <- lm(AveragScore ~ EconomicQuality, data = data_selected)
```

```
# Summarize simple linear regression models
summary(simple_model1)
```

```
##
## Call:
## lm(formula = AveragScore ~ LivingConditions, data = data_selected)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -19.4522  -4.1313  -0.0704   4.0149  10.8746
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   15.98341    1.70103   9.396  <2e-16 ***
## LivingConditions  0.60273    0.02345  25.699  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.969 on 165 degrees of freedom
## Multiple R-squared:  0.8001, Adjusted R-squared:  0.7989
## F-statistic: 660.5 on 1 and 165 DF, p-value: < 2.2e-16
```

```
summary(simple_model2)
```

```
##
## Call:
## lm(formula = AveragScore ~ Health, data = data_selected)
##
## Residuals:
```

```
##      Min      1Q   Median      3Q      Max
## -19.7502 -4.2939 -0.2367   5.1093  15.0967
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -11.09108    3.44662  -3.218  0.00155 **
## Health      1.00761     0.04958  20.325 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.132 on 165 degrees of freedom
## Multiple R-squared:  0.7146, Adjusted R-squared:  0.7128
## F-statistic: 413.1 on 1 and 165 DF,  p-value: < 2.2e-16
```

```
summary(simple_model3)
```

```
##
## Call:
## lm(formula = AveragScore ~ Education, data = data_selected)
##
## Residuals:
##      Min      1Q   Median      3Q      Max
## -15.9192 -3.4137  0.6668   4.1312  11.5349
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  21.9019    1.3497   16.23 <2e-16 ***
## Education     0.6157    0.0218   28.24 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.527 on 165 degrees of freedom
## Multiple R-squared:  0.8286, Adjusted R-squared:  0.8276
## F-statistic: 797.7 on 1 and 165 DF,  p-value: < 2.2e-16
```

```
summary(simple_model4)
```

```
##
## Call:
## lm(formula = AveragScore ~ EconomicQuality, data = data_selected)
##
## Residuals:
##      Min      1Q   Median      3Q      Max
## -15.1391 -4.4007 -0.0824   3.8885  12.2771
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.73001    1.73763   6.751 2.39e-10 ***
## EconomicQuality 0.89834    0.03261  27.544 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.642 on 165 degrees of freedom
```

```
## Multiple R-squared:  0.8214, Adjusted R-squared:  0.8203
## F-statistic: 758.7 on 1 and 165 DF,  p-value: < 2.2e-16
```

```
# Create a table for simple linear regression results
simple_results <- data.frame(
  Model = c("LivingConditions", "Health", "Education", "EconomicQuality"),
  Adj_R2 = c(summary(simple_model1)$adj.r.squared, summary(simple_model2)$adj.r.squared,
             summary(simple_model3)$adj.r.squared, summary(simple_model4)$adj.r.squared),
  Coefficient = c(coef(simple_model1)[2], coef(simple_model2)[2], coef(simple_model3)[2], coef(simple_model4)[2]),
  P_Value = c(summary(simple_model1)$coefficients[2, 4], summary(simple_model2)$coefficients[2, 4],
             summary(simple_model3)$coefficients[2, 4], summary(simple_model4)$coefficients[2, 4])
)

simple_results %>%
  mutate(across(where(is.numeric), round, 4)) %>%
  kable(caption = "Simple Linear Regression Results") %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed", "responsive"))
```

```
## Warning: There was 1 warning in `mutate()`.
## i In argument: `across(where(is.numeric), round, 4)`.
```

## Caused by warning:

```
## ! The `...` argument of `across()` is deprecated as of dplyr 1.1.0.
## Supply arguments directly to `.fns` through an anonymous function instead.
##
## # Previously
##   across(a:b, mean, na.rm = TRUE)
##
## # Now
##   across(a:b, \(x) mean(x, na.rm = TRUE))
```

Table 2: Simple Linear Regression Results

	Model	Adj_R2	Coefficient	P_Value
LivingConditions	LivingConditions	0.7989	0.6027	0
Health	Health	0.7128	1.0076	0
Education	Education	0.8276	0.6157	0
EconomicQuality	EconomicQuality	0.8203	0.8983	0

## Interpretation of the Table

### 1. Model: LivingConditions

- **Adjusted R<sup>2</sup> (0.7989):** This model explains approximately 79.89% of the variance in AverageScore, indicating a strong fit.
- **Coefficient (0.6027):** For every one unit increase in LivingConditions, the AverageScore is expected to increase by 0.6027 units, holding all else constant.
- **P-Value (0):** The relationship between LivingConditions and AverageScore is statistically significant.

### 2. Model: Health

- **Adjusted R<sup>2</sup> (0.7128)**: This model explains approximately 71.28% of the variance in AverageScore, indicating a strong fit but lower than LivingConditions.
- **Coefficient (1.0076)**: For every one unit increase in Health, the AverageScore is expected to increase by 1.0076 units, holding all else constant. This is the highest coefficient among the four models.
- **P-Value (0)**: The relationship between Health and AverageScore is statistically significant.

### 3. Model: Education

- **Adjusted R<sup>2</sup> (0.8276)**: This model explains approximately 82.76% of the variance in AverageScore, the highest among the four models, indicating an excellent fit.
- **Coefficient (0.6157)**: For every one unit increase in Education, the AverageScore is expected to increase by 0.6157 units, holding all else constant.
- **P-Value (0)**: The relationship between Education and AverageScore is statistically significant.

### 4. Model: EconomicQuality

- **Adjusted R<sup>2</sup> (0.8203)**: This model explains approximately 82.03% of the variance in AverageScore, very close to the Education model, indicating an excellent fit.
- **Coefficient (0.8983)**: For every one unit increase in EconomicQuality, the AverageScore is expected to increase by 0.8983 units, holding all else constant.
- **P-Value (0)**: The relationship between EconomicQuality and AverageScore is statistically significant.

## Summary

- **Best Fit**: The Education model has the highest adjusted R-squared value (0.8276), indicating it explains the most variance in AverageScore among the four models.
- **Strong Predictors**: Both Education and EconomicQuality are strong predictors of AverageScore, with adjusted R-squared values above 0.82.
- **Statistical Significance**: All four models show statistically significant relationships with AverageScore, as indicated by their p-values of 0.

## Decision for Multiple Regression

Based on the results of the simple linear regression models, we see that Education, EconomicQuality, LivingConditions, and Health are all significant predictors of AverageScore.

Therefore, we will use these variables as predictors in the multiple regression models. By combining these variables, we aim to develop a more comprehensive model that explains the AverageScore using multiple predictors simultaneously, which may provide better insight and predictive power compared to individual predictors alone. This multiple regression analysis will help us understand the combined effect of these variables and potentially identify the best model based on various criteria such as Adjusted R<sup>2</sup>, AIC, BIC, Mallow's Cp, and VIF.

```
# Multiple linear regression model
multiple_model <- lm(AverageScore ~ LivingConditions + Health + Education + EconomicQuality, data = data)
summary(multiple_model)
```

```
##
## Call:
## lm(formula = AveragScore ~ LivingConditions + Health + Education +
##     EconomicQuality, data = data_selected)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.8373  -2.6251   0.3614   3.0011   7.8887
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    8.72104    2.72752   3.197 0.001668 **
## LivingConditions 0.12701    0.05313   2.391 0.017972 *
## Health          0.08006    0.06933   1.155 0.249893
## Education       0.19650    0.05186   3.789 0.000213 ***
## EconomicQuality 0.45446    0.04428  10.264 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.145 on 162 degrees of freedom
## Multiple R-squared:  0.9054, Adjusted R-squared:  0.903
## F-statistic: 387.5 on 4 and 162 DF,  p-value: < 2.2e-16
```

## Explanation of Multiple Linear Regression Results

The multiple linear regression model includes four predictors: LivingConditions, Health, Education, and EconomicQuality. The model aims to predict the AverageScore based on these predictors. Below are the detailed results and interpretation:

### Model Summary

- **Residuals:**
  - **Min:** -10.8373
  - **1Q (First Quartile):** -2.6251
  - **Median:** 0.3614
  - **3Q (Third Quartile):** 3.0011
  - **Max:** 7.8887

### Coefficients

- **Intercept:** 8.72104
  - **Std. Error:** 2.72752
  - **t value:** 3.197
  - **Pr(>|t|):** 0.001668 \*\* (significant at 0.01 level)
- **LivingConditions:** 0.12701
  - **Std. Error:** 0.05313

- **t value:** 2.391
- **Pr(>|t|):** 0.017972 \* (significant at 0.05 level)
- **Health:** 0.08006
  - **Std. Error:** 0.06933
  - **t value:** 1.155
  - **Pr(>|t|):** 0.249893 (not significant)
- **Education:** 0.19650
  - **Std. Error:** 0.05186
  - **t value:** 3.789
  - **Pr(>|t|):** 0.000213 \*\*\* (significant at 0.001 level)
- **EconomicQuality:** 0.45446
  - **Std. Error:** 0.04428
  - **t value:** 10.264
  - **Pr(>|t|):** < 2e-16 \*\*\* (highly significant)

### Model Fit

- **Residual standard error:** 4.145 on 162 degrees of freedom
- **Multiple R-squared:** 0.9054
- **Adjusted R-squared:** 0.903
- **F-statistic:** 387.5 on 4 and 162 DF
- **p-value:** < 2.2e-16

### Interpretation

- **Model Significance:** The overall model is highly significant with a p-value < 2.2e-16, indicating that the predictors collectively explain a significant portion of the variance in AverageScore.
- **Adjusted R-squared (0.903):** This indicates that approximately 90.3% of the variance in AverageScore is explained by the model, which is a very strong fit.
- **Significant Predictors:**
  - **Education** and **EconomicQuality** are highly significant predictors with p-values less than 0.001.
  - **LivingConditions** is also significant at the 0.05 level.
  - **Health** is not a significant predictor with a p-value of 0.249893.
- **Coefficients:**
  - **EconomicQuality** has the highest coefficient (0.45446), indicating it has the strongest impact on AverageScore among the predictors.
  - **Education** also has a notable positive impact (0.19650).
  - **LivingConditions** and **Health** have smaller coefficients, with Health not being statistically significant.

## Simplification of the Model

Based on these results, **Health** can be considered for removal from the model to simplify it, as it is not a significant predictor (p-value = 0.249893).

The revised model will exclude the Health variable and re-evaluate the model fit and significance:

**Simplified Model Formula** By removing the non-significant predictor (Health), we aim to simplify the model while retaining the significant predictors that contribute meaningfully to explaining the variance in AverageScore.

```
# Simplified multiple linear regression model (removing 'Health')
simplified_model <- lm(AveragScore ~ LivingConditions + Education + EconomicQuality, data = data_selected)
summary(simplified_model)
```

```
##
## Call:
## lm(formula = AveragScore ~ LivingConditions + Education + EconomicQuality,
##     data = data_selected)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11.1462  -2.8974   0.2246   3.1328   7.7585
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    11.42128     1.40547   8.126 1.05e-13 ***
## LivingConditions  0.15459     0.04751   3.254  0.00138 **
## Education        0.20958     0.05066   4.137 5.63e-05 ***
## EconomicQuality  0.45642     0.04429  10.305 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.149 on 163 degrees of freedom
## Multiple R-squared:  0.9046, Adjusted R-squared:  0.9028
## F-statistic: 515.1 on 3 and 163 DF, p-value: < 2.2e-16
```

```
# Calculate VIF for the simplified model
vif_simplified_model <- vif(simplified_model)
```

## Model Summary

- Residuals:
  - Min: -11.1462
  - 1Q (First Quartile): -2.8974
  - Median: 0.2246
  - 3Q (Third Quartile): 3.1328
  - Max: 7.7585



## Coefficients

- **Intercept:**
  - **Estimate:** 11.42128
  - **Std. Error:** 1.40547
  - **t value:** 8.126
  - **Pr(>|t|):** 1.05e-13 (highly significant)
- **LivingConditions:**
  - **Estimate:** 0.15459
  - **Std. Error:** 0.04751
  - **t value:** 3.254
  - **Pr(>|t|):** 0.00138 (significant)
- **Education:**
  - **Estimate:** 0.20958
  - **Std. Error:** 0.05066
  - **t value:** 4.137
  - **Pr(>|t|):** 5.63e-05 (highly significant)
- **EconomicQuality:**
  - **Estimate:** 0.45642
  - **Std. Error:** 0.04429
  - **t value:** 10.305
  - **Pr(>|t|):** < 2e-16 (highly significant)

## Model Fit

- **Residual standard error:** 4.149 on 163 degrees of freedom
- **Multiple R-squared:** 0.9046
- **Adjusted R-squared:** 0.9028
- **F-statistic:** 515.1 on 3 and 163 DF
- **p-value:** < 2.2e-16 (model is highly significant)

## Interpretation

- **Intercept:** The intercept value is 11.42128, meaning that when all predictor variables are zero, the average score is 11.42128.
- **LivingConditions:** The coefficient for LivingConditions is 0.15459, indicating that a one-unit increase in LivingConditions is associated with a 0.15459 increase in the average score, holding all other variables constant. The p-value of 0.00138 indicates that this relationship is statistically significant.

- **Education:** The coefficient for Education is 0.20958, suggesting that a one-unit increase in Education is associated with a 0.20958 increase in the average score, holding all other variables constant. The p-value of 5.63e-05 indicates that this relationship is highly significant.
- **EconomicQuality:** The coefficient for EconomicQuality is 0.45642, indicating that a one-unit increase in EconomicQuality is associated with a 0.45642 increase in the average score, holding all other variables constant. The p-value of  $< 2e-16$  indicates that this relationship is highly significant.
- **Model Fit:** The Adjusted R-squared value of 0.9028 means that approximately 90.28% of the variability in the average score is explained by the model. The F-statistic of 515.1 with a p-value  $< 2.2e-16$  shows that the model is highly significant.

## Conclusion

The simplified multiple linear regression model shows that LivingConditions, Education, and EconomicQuality are significant predictors of the average score. Among them, EconomicQuality has the largest effect size. Based on this analysis, we have removed the Health variable from the previous model as it was not a significant predictor.

```
# Define the full model formula dynamically
predictors <- setdiff(names(data_selected), "AveragScore")
full_model_formula <- as.formula(paste("AveragScore ~", paste(predictors, collapse = " + ")))

null_model <- lm(AveragScore ~ 1, data = data_selected)
full_model <- lm(AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
  SafetySecurity + SocialCapital + MarketAccessInfrastructure +
  Governance + EconomicQuality + Health + NaturalEnvironment +
  LivingConditions + EnterpriseConditions, data = data_selected)

# Perform forward stepwise regression using AIC
forward_stepwise_aic <- stepAIC(null_model, scope = list(lower = null_model, upper = full_model), direc
```

```
## Start:  AIC=865.56
## AveragScore ~ 1
##
##
##      Df Sum of Sq    RSS   AIC
## + InvestmentEnvironment      1    26465 2942.7 483.14
## + MarketAccessInfrastructure  1    26170 3238.0 499.11
## + Governance                  1    25111 4296.9 546.36
## + Education                   1    24367 5040.4 573.01
## + EconomicQuality             1    24155 5253.1 579.91
## + LivingConditions            1    23529 5878.3 598.69
## + EnterpriseConditions        1    22983 6424.9 613.54
## + Health                      1    21014 8393.6 658.18
## + SafetySecurity              1    20191 9216.8 673.80
## + SocialCapital               1    17192 12215.8 720.85
## + PersonelFreedom             1    17125 12283.3 721.77
## + NaturalEnvironment          1    14077 15330.6 758.78
## <none>                        29407.8 865.56
##
## Step:  AIC=483.14
## AveragScore ~ InvestmentEnvironment
##
```

```

##                                Df Sum of Sq    RSS    AIC
## + Education                   1   1269.05 1673.6 390.89
## + SafetySecurity              1   1191.74 1750.9 398.43
## + LivingConditions            1   1142.21 1800.5 403.09
## + SocialCapital               1    945.19 1997.5 420.44
## + PersonelFreedom            1    866.96 2075.7 426.85
## + MarketAccessInfrastructure  1    817.35 2125.3 430.79
## + NaturalEnvironment          1    772.65 2170.0 434.27
## + Health                     1    749.28 2193.4 436.06
## + Governance                 1    660.38 2282.3 442.69
## + EconomicQuality            1    539.86 2402.8 451.29
## <none>                        2942.7 483.14
## + EnterpriseConditions        1      1.15 2941.5 485.07
##
## Step:   AIC=390.89
## AveragScore ~ InvestmentEnvironment + Education
##
##                                Df Sum of Sq    RSS    AIC
## + PersonelFreedom            1    964.83  708.78 249.41
## + Governance                 1    833.17  840.44 277.86
## + SafetySecurity             1    695.29  978.32 303.23
## + NaturalEnvironment         1    609.03 1064.58 317.34
## + SocialCapital              1    520.64 1152.97 330.66
## + EconomicQuality            1    127.07 1546.54 379.71
## + EnterpriseConditions        1    119.65 1553.96 380.50
## + LivingConditions           1     79.74 1593.87 384.74
## + MarketAccessInfrastructure  1     42.93 1630.68 388.55
## + Health                     1     24.41 1649.20 390.44
## <none>                        1673.61 390.89
##
## Step:   AIC=249.41
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom
##
##                                Df Sum of Sq    RSS    AIC
## + SafetySecurity             1    272.221 436.56 170.48
## + EconomicQuality            1    179.952 528.83 202.50
## + SocialCapital              1    165.922 542.86 206.87
## + Health                     1    120.927 587.86 220.17
## + Governance                 1    119.181 589.60 220.66
## + MarketAccessInfrastructure  1    118.768 590.02 220.78
## + NaturalEnvironment         1     81.741 627.04 230.94
## + LivingConditions           1     74.926 633.86 232.75
## + EnterpriseConditions        1     67.048 641.74 234.81
## <none>                        708.78 249.41
##
## Step:   AIC=170.48
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##           SafetySecurity
##
##                                Df Sum of Sq    RSS    AIC
## + SocialCapital              1    134.260 302.30 111.10
## + EconomicQuality            1    121.511 315.05 118.00
## + MarketAccessInfrastructure  1    117.259 319.30 120.24
## + Health                     1     99.255 337.31 129.40

```

```

## + NaturalEnvironment      1    61.145 375.42 147.28
## + Governance              1    55.365 381.20 149.83
## + LivingConditions         1    54.874 381.69 150.04
## + EnterpriseConditions     1    41.317 395.25 155.87
## <none>                    436.56 170.48
##
## Step: AIC=111.1
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital
##
##              Df Sum of Sq    RSS    AIC
## + MarketAccessInfrastructure  1   112.974 189.33  34.956
## + EconomicQuality            1    73.051 229.25  66.910
## + Health                     1    72.015 230.29  67.663
## + LivingConditions            1    57.056 245.25  78.173
## + EnterpriseConditions        1    41.164 261.14  88.658
## + Governance                  1    36.458 265.84  91.641
## + NaturalEnvironment          1    33.574 268.73  93.443
## <none>                      302.30 111.103
##
## Step: AIC=34.96
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure
##
##              Df Sum of Sq    RSS    AIC
## + Governance                  1    55.825 133.50 -21.386
## + EconomicQuality             1    46.814 142.51 -10.479
## + EnterpriseConditions         1    31.525 157.80   6.539
## + Health                     1    29.803 159.53   8.352
## + NaturalEnvironment          1    22.087 167.24  16.240
## + LivingConditions            1    11.588 177.74  26.409
## <none>                      189.33  34.956
##
## Step: AIC=-21.39
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance
##
##              Df Sum of Sq    RSS    AIC
## + EconomicQuality             1    32.412 101.09 -65.830
## + Health                     1    31.688 101.81 -64.639
## + LivingConditions            1    27.734 105.77 -58.274
## + NaturalEnvironment          1    23.111 110.39 -51.130
## + EnterpriseConditions        1     3.960 129.54 -24.415
## <none>                      133.50 -21.386
##
## Step: AIC=-65.83
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality
##
##              Df Sum of Sq    RSS    AIC
## + Health                     1    40.831  60.259 -150.230
## + LivingConditions            1    28.408  72.683 -118.926

```

```

## + NaturalEnvironment      1      19.767  81.323 -100.167
## + EnterpriseConditions    1       2.637  98.454 -68.243
## <none>                    101.091  -65.830
##
## Step:  AIC=-150.23
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##      SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##      Governance + EconomicQuality + Health
##
##              Df Sum of Sq    RSS    AIC
## + NaturalEnvironment      1   24.0684 36.191 -233.37
## + LivingConditions         1    9.8852 50.374 -178.15
## + EnterpriseConditions     1    6.9876 53.272 -168.81
## <none>                    60.259 -150.23
##
## Step:  AIC=-233.37
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##      SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##      Governance + EconomicQuality + Health + NaturalEnvironment
##
##              Df Sum of Sq    RSS    AIC
## + LivingConditions         1   20.8682 15.323 -374.91
## + EnterpriseConditions     1    6.3207 29.870 -263.43
## <none>                    36.191 -233.37
##
## Step:  AIC=-374.91
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##      SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##      Governance + EconomicQuality + Health + NaturalEnvironment +
##      LivingConditions
##
##              Df Sum of Sq    RSS    AIC
## + EnterpriseConditions     1    15.321  0.0013 -1935.02
## <none>                    15.3226  -374.91
##
## Step:  AIC=-1935.02
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##      SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##      Governance + EconomicQuality + Health + NaturalEnvironment +
##      LivingConditions + EnterpriseConditions

```

## Stepwise Regression Explanation

Stepwise regression is a method of fitting regression models in which the choice of predictive variables is carried out by an automatic procedure. In this case, it starts with no predictors and adds them one at a time (forward selection) based on the Akaike Information Criterion (AIC), a measure of the quality of a statistical model.

Here are the steps taken by the forward stepwise regression:

### 1. Initial Model (Intercept Only)

- AIC = 865.56
- Starting with only the intercept, no predictors.

## 2. Step 1: Add InvestmentEnvironment

- Adding InvestmentEnvironment reduces the AIC from 865.56 to 483.14.
- This is the predictor that provides the largest reduction in AIC.

## 3. Step 2: Add Education

- Adding Education to the model that already includes InvestmentEnvironment further reduces the AIC to 390.89.
- This step adds the predictor that provides the next largest reduction in AIC.

## 4. Step 3: Add PersonelFreedom

- Adding PersonelFreedom to the model that already includes InvestmentEnvironment and Education reduces the AIC to 249.41.

## 5. Step 4: Add SafetySecurity

- Adding SafetySecurity to the model reduces the AIC to 170.48.

## 6. Step 5: Add SocialCapital

- Adding SocialCapital reduces the AIC to 111.10.

## 7. Step 6: Add MarketAccessInfrastructure

- Adding MarketAccessInfrastructure reduces the AIC to 34.96.

## 8. Step 7: Add Governance

- Adding Governance reduces the AIC to -21.39.

## 9. Step 8: Add EconomicQuality

- Adding EconomicQuality reduces the AIC to -65.83.

## 10. Step 9: Add Health

- Adding Health reduces the AIC to -150.23.

## 11. Step 10: Add NaturalEnvironment

- Adding NaturalEnvironment reduces the AIC to -233.37.

## 12. Step 11: Add LivingConditions

- Adding LivingConditions reduces the AIC to -374.91.

## 13. Step 12: Add EnterpriseConditions

- Adding EnterpriseConditions reduces the AIC to -1935.02.

The final model includes all the predictors listed, with an AIC of -1935.02. This indicates a very good fit of the model to the data.

## Interpretation of the Final Model

The final model is:  $\text{AveragScore} = \text{InvestmentEnvironment} + \text{Education} + \text{PersonelFreedom} + \text{SafetySecurity} + \text{SocialCapital} + \text{MarketAccessInfrastructure} + \text{Governance} + \text{EconomicQuality} + \text{Health} + \text{NaturalEnvironment} + \text{LivingConditions} + \text{EnterpriseConditions}$

**Coefficients:** Each predictor has a coefficient, which indicates the change in the response variable for a one-unit change in the predictor, holding all other predictors constant.

For example:

- **InvestmentEnvironment:** Adding this predictor in the initial step decreased AIC significantly, indicating it is a strong predictor of **AveragScore**.
- **Education:** Also shows significant influence on **AveragScore** when added after **InvestmentEnvironment**.

**AIC:** AIC is a measure of the relative quality of statistical models for a given set of data. Lower AIC indicates a better model fit. The stepwise selection process chooses predictors to minimize the AIC, thus improving model quality at each step.

## Summary

The stepwise regression shows the process of building a model by adding one predictor at a time based on the AIC criterion, leading to a final model that includes all 12 predictors. Each step's selection is justified by the reduction in AIC, showing the model improvement.

The final model with the lowest AIC includes:

- InvestmentEnvironment
- Education
- PersonelFreedom
- SafetySecurity
- SocialCapital
- MarketAccessInfrastructure
- Governance
- EconomicQuality
- Health
- NaturalEnvironment
- LivingConditions
- EnterpriseConditions

These predictors together provide the best fit for predicting **AveragScore**, considering the AIC criterion.

This stepwise process ensures that each predictor added significantly improves the model fit, justifying its inclusion. The final model has a significantly lower AIC compared to the intercept-only model, indicating a well-fitted model.

```
# Perform backward stepwise regression using AIC
backward_stepwise_aic <- stepAIC(full_model, direction = "backward")
```

```
## Start:  AIC=-1935.02
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##      SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##      Governance + EconomicQuality + Health + NaturalEnvironment +
##      LivingConditions + EnterpriseConditions
##
##              Df Sum of Sq      RSS      AIC
## <none>              0.001 -1935.02
## - MarketAccessInfrastructure  1    14.841  14.842  -380.23
## - EnterpriseConditions        1    15.321  15.323  -374.91
## - InvestmentEnvironment       1    16.921  16.923  -358.32
## - Health                     1    21.900  21.901  -315.25
## - Governance                 1    22.802  22.804  -308.51
## - LivingConditions           1    29.869  29.870  -263.43
## - EconomicQuality            1    31.882  31.883  -252.54
## - NaturalEnvironment         1    37.179  37.180  -226.87
## - Education                  1    40.898  40.899  -210.95
## - SocialCapital              1    49.957  49.959  -177.53
## - PersonelFreedom            1    88.682  88.683   -81.70
## - SafetySecurity             1   124.461 124.463   -25.10
```

### Initial Model (Full Model)

The initial model includes all predictors: AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity

- **AIC = -1935.02**

### Steps in Backward Elimination

**Step 1: Evaluate Removal of Each Predictor** Each row in the output shows the result of removing one predictor at a time from the full model. The AIC for the resulting model is calculated if that predictor is removed.

- **<none>**: Indicates the AIC of the current model with all predictors.
- Each line below “<none>” indicates the AIC of the model if the corresponding predictor were removed.

For example:

- Removing **MarketAccessInfrastructure** results in an AIC of -380.23.
- Removing **EnterpriseConditions** results in an AIC of -374.91.
- Removing **InvestmentEnvironment** results in an AIC of -358.32.

**Interpretation of Each Step** The process evaluates the impact of removing each predictor:

- **MarketAccessInfrastructure**: Removing this predictor results in the smallest increase in AIC to -380.23, making it a strong candidate for removal.
- **EnterpriseConditions**: Removing this predictor increases the AIC slightly more to -374.91.
- **InvestmentEnvironment**: Removing this predictor results in an AIC of -358.32, showing a more significant impact compared to the previous ones.



## Key Points in Backward Elimination

1. **Start with All Predictors:** The initial model starts with all available predictors, having the lowest AIC of -1935.02.
2. **Evaluate Removal Impact:** Each predictor is evaluated for removal, calculating the resulting AIC. The predictor whose removal causes the smallest increase in AIC is considered for elimination.
3. **Continue Until No Improvement:** The process continues until removing any further predictor would result in a higher AIC, indicating no improvement.
4. **Final Decision:**
  - The model with the lowest AIC (-1935.02) is chosen. This model includes all predictors as none of the removals resulted in a lower AIC.

## Conclusion

In this backward elimination stepwise regression:

- The final model includes all the predictors since none of the removals resulted in a better (lower) AIC than the full model.
- This indicates that each predictor contributes significantly to the model, and removing any predictor would reduce the model's quality as indicated by an increase in AIC.

## Summary of Final Model Predictors

The final model retains all predictors:

- InvestmentEnvironment
- Education
- PersonelFreedom
- SafetySecurity
- SocialCapital
- MarketAccessInfrastructure
- Governance
- EconomicQuality
- Health
- NaturalEnvironment
- LivingConditions
- EnterpriseConditions

This comprehensive model is considered the best fit given the data, as removing any predictor would lead to a higher AIC and thus a worse model fit.

```
# Perform both-direction stepwise regression using AIC
```

```
both_stepwise_aic <- stepAIC(null_model, scope = list(lower = null_model, upper = full_model), direction
```

```
## Start: AIC=865.56
```

```
## AveragScore ~ 1
```

```
##
```

	Df	Sum of Sq	RSS	AIC
## + InvestmentEnvironment	1	26465	2942.7	483.14
## + MarketAccessInfrastructure	1	26170	3238.0	499.11
## + Governance	1	25111	4296.9	546.36
## + Education	1	24367	5040.4	573.01
## + EconomicQuality	1	24155	5253.1	579.91
## + LivingConditions	1	23529	5878.3	598.69
## + EnterpriseConditions	1	22983	6424.9	613.54
## + Health	1	21014	8393.6	658.18
## + SafetySecurity	1	20191	9216.8	673.80
## + SocialCapital	1	17192	12215.8	720.85
## + PersonelFreedom	1	17125	12283.3	721.77
## + NaturalEnvironment	1	14077	15330.6	758.78
## <none>			29407.8	865.56

```
##
```

```
## Step: AIC=483.14
```

```
## AveragScore ~ InvestmentEnvironment
```

```
##
```

	Df	Sum of Sq	RSS	AIC
## + Education	1	1269.1	1673.6	390.89
## + SafetySecurity	1	1191.7	1750.9	398.43
## + LivingConditions	1	1142.2	1800.5	403.09
## + SocialCapital	1	945.2	1997.5	420.44
## + PersonelFreedom	1	867.0	2075.7	426.85
## + MarketAccessInfrastructure	1	817.4	2125.3	430.79
## + NaturalEnvironment	1	772.6	2170.0	434.27
## + Health	1	749.3	2193.4	436.06
## + Governance	1	660.4	2282.3	442.69
## + EconomicQuality	1	539.9	2402.8	451.29
## <none>			2942.7	483.14
## + EnterpriseConditions	1	1.1	2941.5	485.07
## - InvestmentEnvironment	1	26465.1	29407.8	865.56

```
##
```

```
## Step: AIC=390.89
```

```
## AveragScore ~ InvestmentEnvironment + Education
```

```
##
```

	Df	Sum of Sq	RSS	AIC
## + PersonelFreedom	1	964.8	708.8	249.41
## + Governance	1	833.2	840.4	277.86
## + SafetySecurity	1	695.3	978.3	303.23
## + NaturalEnvironment	1	609.0	1064.6	317.34
## + SocialCapital	1	520.6	1153.0	330.66
## + EconomicQuality	1	127.1	1546.5	379.71
## + EnterpriseConditions	1	119.7	1554.0	380.50
## + LivingConditions	1	79.7	1593.9	384.74
## + MarketAccessInfrastructure	1	42.9	1630.7	388.55
## + Health	1	24.4	1649.2	390.44

```

## <none>                                1673.6 390.89
## - Education                           1    1269.1 2942.7 483.14
## - InvestmentEnvironment                1    3366.8 5040.4 573.01
##
## Step: AIC=249.41
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom
##
##              Df Sum of Sq    RSS    AIC
## + SafetySecurity      1    272.22  436.56 170.48
## + EconomicQuality      1    179.95  528.83 202.50
## + SocialCapital        1    165.92  542.86 206.87
## + Health               1    120.93  587.86 220.17
## + Governance           1    119.18  589.60 220.66
## + MarketAccessInfrastructure 1    118.77  590.02 220.78
## + NaturalEnvironment    1     81.74  627.04 230.94
## + LivingConditions       1     74.93  633.86 232.75
## + EnterpriseConditions   1     67.05  641.74 234.81
## <none>                                708.78 249.41
## - PersonelFreedom        1    964.83 1673.61 390.89
## - InvestmentEnvironment    1   1360.55 2069.34 426.34
## - Education               1   1366.92 2075.70 426.85
##
## Step: AIC=170.48
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity
##
##              Df Sum of Sq    RSS    AIC
## + SocialCapital        1    134.26  302.30 111.10
## + EconomicQuality      1    121.51  315.05 118.00
## + MarketAccessInfrastructure 1    117.26  319.30 120.24
## + Health               1     99.26  337.31 129.40
## + NaturalEnvironment    1     61.14  375.42 147.28
## + Governance           1     55.36  381.20 149.83
## + LivingConditions       1     54.87  381.69 150.04
## + EnterpriseConditions   1     41.32  395.25 155.87
## <none>                                436.56 170.48
## - SafetySecurity        1    272.22  708.78 249.41
## - PersonelFreedom        1    541.76  978.32 303.23
## - Education              1    947.10 1383.66 361.12
## - InvestmentEnvironment    1   1116.64 1553.20 380.42
##
## Step: AIC=111.1
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital
##
##              Df Sum of Sq    RSS    AIC
## + MarketAccessInfrastructure 1    112.97  189.33  34.96
## + EconomicQuality            1     73.05  229.25  66.91
## + Health                     1     72.01  230.29  67.66
## + LivingConditions            1     57.06  245.25  78.17
## + EnterpriseConditions        1     41.16  261.14  88.66
## + Governance                  1     36.46  265.84  91.64
## + NaturalEnvironment          1     33.57  268.73  93.44
## <none>                        302.30 111.10

```

```

## - SocialCapital          1    134.26  436.56 170.48
## - SafetySecurity         1    240.56  542.86 206.87
## - PersonelFreedom        1    349.04  651.34 237.29
## - Education              1    733.73 1036.03 314.80
## - InvestmentEnvironment   1   1045.56 1347.86 358.74
##
## Step: AIC=34.96
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure
##
##              Df Sum of Sq    RSS    AIC
## + Governance      1     55.82 133.50 -21.386
## + EconomicQuality  1     46.81 142.51 -10.479
## + EnterpriseConditions  1     31.53 157.80   6.539
## + Health           1     29.80 159.52   8.352
## + NaturalEnvironment  1     22.09 167.24  16.240
## + LivingConditions   1     11.59 177.74  26.409
## <none>                189.33  34.956
## - MarketAccessInfrastructure  1    112.97 302.30 111.103
## - SocialCapital        1    129.98 319.30 120.241
## - Education            1    199.51 388.84 153.142
## - InvestmentEnvironment  1    232.71 422.04 166.827
## - SafetySecurity        1    239.63 428.96 169.544
## - PersonelFreedom       1    395.81 585.14 221.394
##
## Step: AIC=-21.39
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance
##
##              Df Sum of Sq    RSS    AIC
## + EconomicQuality      1     32.412 101.09 -65.830
## + Health                 1     31.688 101.81 -64.639
## + LivingConditions       1     27.734 105.77 -58.274
## + NaturalEnvironment     1     23.111 110.39 -51.130
## + EnterpriseConditions    1      3.960 129.54 -24.415
## <none>                  133.50 -21.386
## - InvestmentEnvironment  1     49.234 182.74  29.039
## - Governance             1     55.825 189.33  34.956
## - SocialCapital          1    107.064 240.57  74.956
## - MarketAccessInfrastructure  1    132.341 265.84  91.641
## - PersonelFreedom        1    145.657 279.16  99.803
## - SafetySecurity          1    182.417 315.92 120.462
## - Education              1    209.008 342.51 133.958
##
## Step: AIC=-65.83
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality
##
##              Df Sum of Sq    RSS    AIC
## + Health                1     40.831  60.259 -150.230
## + LivingConditions       1     28.408  72.683 -118.926
## + NaturalEnvironment     1     19.767  81.323 -100.167

```

```

## + EnterpriseConditions      1      2.637  98.454 -68.243
## <none>                      1      101.091 -65.830
## - EconomicQuality          1     32.412 133.503 -21.386
## - InvestmentEnvironment     1     33.156 134.246 -20.459
## - Governance                1     41.423 142.513 -10.479
## - SocialCapital             1     78.712 179.802  28.335
## - MarketAccessInfrastructure 1    104.134 205.224  50.420
## - SafetySecurity            1    165.105 266.196  93.862
## - PersonelFreedom           1    170.347 271.437  97.118
## - Education                 1    189.546 290.636 108.531
##
## Step: AIC=-150.23
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health
##
##              Df Sum of Sq      RSS      AIC
## + NaturalEnvironment      1    24.068  36.191 -233.374
## + LivingConditions         1     9.885  50.374 -178.153
## + EnterpriseConditions      1     6.988  53.272 -168.813
## <none>                     1    60.259 -150.230
## - InvestmentEnvironment     1    33.967  94.226 -77.573
## - Health                   1    40.831 101.091 -65.830
## - EconomicQuality          1    41.555 101.814 -64.639
## - Governance               1    41.570 101.829 -64.614
## - MarketAccessInfrastructure 1    57.671 117.930 -40.100
## - SocialCapital            1    59.684 119.943 -37.273
## - Education                1    98.818 159.077   9.883
## - SafetySecurity           1   152.700 212.959  58.599
## - PersonelFreedom          1   195.424 255.684  89.133
##
## Step: AIC=-233.37
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health + NaturalEnvironment
##
##              Df Sum of Sq      RSS      AIC
## + LivingConditions         1    20.868  15.323 -374.91
## + EnterpriseConditions      1     6.321  29.870 -263.43
## <none>                     1    36.191 -233.37
## - NaturalEnvironment      1    24.068  60.259 -150.23
## - InvestmentEnvironment     1    37.308  73.498 -117.06
## - EconomicQuality          1    37.961  74.152 -115.58
## - Governance               1    43.144  79.335 -104.30
## - Health                   1    45.132  81.323 -100.17
## - SocialCapital            1    45.566  81.757  -99.28
## - MarketAccessInfrastructure 1    49.385  85.576  -91.65
## - Education                1    98.190 134.381 -16.29
## - PersonelFreedom          1   118.772 154.963   7.51
## - SafetySecurity           1   144.796 180.987  33.43
##
## Step: AIC=-374.91
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +

```

```

##      Governance + EconomicQuality + Health + NaturalEnvironment +
##      LivingConditions
##
##
##              Df Sum of Sq      RSS      AIC
## + EnterpriseConditions      1    15.321    0.001 -1935.02
## <none>                      15.323    -374.91
## - Health                    1    20.420   35.743  -235.46
## - LivingConditions          1    20.868   36.191  -233.37
## - MarketAccessInfrastructure 1    26.222   41.545  -210.33
## - EconomicQuality           1    34.795   50.118  -179.01
## - NaturalEnvironment        1    35.051   50.374  -178.15
## - InvestmentEnvironment     1    38.540   53.862  -166.97
## - Education                 1    38.876   54.199  -165.93
## - SocialCapital             1    46.136   61.459  -144.94
## - Governance                1    56.541   71.863  -118.82
## - PersonelFreedom           1    78.035   93.357   -75.12
## - SafetySecurity            1   127.607  142.930    -3.99
##
## Step:  AIC=-1935.02
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##      SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##      Governance + EconomicQuality + Health + NaturalEnvironment +
##      LivingConditions + EnterpriseConditions
##
##              Df Sum of Sq      RSS      AIC
## <none>                      0.001 -1935.02
## - MarketAccessInfrastructure 1    14.841   14.842  -380.23
## - EnterpriseConditions      1    15.321   15.323  -374.91
## - InvestmentEnvironment     1    16.921   16.923  -358.32
## - Health                    1    21.900   21.901  -315.25
## - Governance                1    22.802   22.804  -308.51
## - LivingConditions          1    29.869   29.870  -263.43
## - EconomicQuality           1    31.882   31.883  -252.54
## - NaturalEnvironment        1    37.179   37.180  -226.87
## - Education                 1    40.898   40.899  -210.95
## - SocialCapital             1    49.957   49.959  -177.53
## - PersonelFreedom           1    88.682   88.683   -81.70
## - SafetySecurity            1   124.461  124.463   -25.10

```

## Both Direction Stepwise Regression Explanation

Both direction stepwise regression combines both forward selection and backward elimination methods to find the best fitting model based on the Akaike Information Criterion (AIC). Here's a detailed explanation of each step:

**Initial Model (Null Model)** The initial model starts with no predictors (null model), and the AIC is 865.56.

AveragScore 1

## Step-by-Step Process

### 1. Step 1 (Forward Selection):

- Evaluate each predictor to see which one, if added to the model, would result in the largest decrease in AIC.
- **InvestmentEnvironment** is added, reducing the AIC to 483.14.

AveragScore InvestmentEnvironment

## 2. Step 2:

- Evaluate the addition of each remaining predictor.
- **Education** is added next, reducing the AIC to 390.89.

AveragScore InvestmentEnvironment+Education

## 3. Step 3:

- Evaluate the addition of each remaining predictor.
- **PersonelFreedom** is added, reducing the AIC to 249.41.

AveragScore InvestmentEnvironment+Education+PersonelFreedom

## 4. Step 4:

- Evaluate the addition of each remaining predictor.
- **SafetySecurity** is added, reducing the AIC to 170.48.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity

## 5. Step 5:

- Evaluate the addition of each remaining predictor.
- **SocialCapital** is added, reducing the AIC to 111.1.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital

## 6. Step 6:

- Evaluate the addition of each remaining predictor.
- **MarketAccessInfrastructure** is added, reducing the AIC to 34.96.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital+MarketAccessInfra

## 7. Step 7:

- Evaluate the addition of each remaining predictor.
- **Governance** is added, reducing the AIC to -21.39.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital+MarketAccessInfra

## 8. Step 8:

- Evaluate the addition of each remaining predictor.
- **EconomicQuality** is added, reducing the AIC to -65.83.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital+MarketAccessInfra

## 9. Step 9:

- Evaluate the addition of each remaining predictor.

- **Health** is added, reducing the AIC to -150.23.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital+MarketAccessInfra

10. **Step 10:**

- Evaluate the addition of each remaining predictor.
- **NaturalEnvironment** is added, reducing the AIC to -233.37.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital+MarketAccessInfra

11. **Step 11:**

- Evaluate the addition of each remaining predictor.
- **LivingConditions** is added, reducing the AIC to -374.91.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital+MarketAccessInfra

12. **Step 12:**

- Evaluate the addition of each remaining predictor.
- **EnterpriseConditions** is added, reducing the AIC to -1935.02.

AveragScore InvestmentEnvironment+Education+PersonelFreedom+SafetySecurity+SocialCapital+MarketAccessInfra

**Evaluation of Removing Each Predictor (Backward Elimination)** After reaching the final model with all predictors included, the stepwise regression evaluates the impact of removing each predictor one by one:

1. **MarketAccessInfrastructure:** Removing this predictor results in an AIC of -380.23.
2. **EnterpriseConditions:** Removing this predictor results in an AIC of -374.91.
3. **InvestmentEnvironment:** Removing this predictor results in an AIC of -358.32.
4. **Health:** Removing this predictor results in an AIC of -315.25.
5. **Governance:** Removing this predictor results in an AIC of -308.51.
6. **LivingConditions:** Removing this predictor results in an AIC of -263.43.
7. **EconomicQuality:** Removing this predictor results in an AIC of -252.54.
8. **NaturalEnvironment:** Removing this predictor results in an AIC of -226.87.
9. **Education:** Removing this predictor results in an AIC of -210.95.
10. **SocialCapital:** Removing this predictor results in an AIC of -177.53.
11. **PersonelFreedom:** Removing this predictor results in an AIC of -81.70.
12. **SafetySecurity:** Removing this predictor results in an AIC of -25.10.



## Summary

- **Both direction stepwise regression** starts with no predictors and adds them one by one to minimize the AIC.
- **Forward Selection:**
  - Predictors added in order: InvestmentEnvironment, Education, PersonelFreedom, SafetySecurity, SocialCapital, MarketAccessInfrastructure, Governance, EconomicQuality, Health, NaturalEnvironment, LivingConditions, EnterpriseConditions.
- **Backward Elimination:**
  - After including all predictors, the model checks if removing any predictor improves the AIC.
  - None of the predictors should be removed as it would increase the AIC, confirming the optimal model found.

## Conclusion

The final model includes all the predictors, indicating each predictor contributes significantly to explaining the variance in the dependent variable, AveragScore. Removing any of these predictors would result in a higher AIC, indicating a poorer model fit.

```
# Fit the final model
final_model <- lm(AveragScore ~ InvestmentEnvironment + Education +
                  PersonelFreedom + SafetySecurity + SocialCapital +
                  MarketAccessInfrastructure + Governance + EconomicQuality +
                  Health + NaturalEnvironment + LivingConditions +
                  EnterpriseConditions, data = data_selected)

# Calculate VIF for each predictor
vif_values <- vif(final_model)
print(vif_values)
```

```
##      InvestmentEnvironment      Education
##      19.095243                10.903521
##      PersonelFreedom          SafetySecurity
##      5.157736                 2.848161
##      SocialCapital MarketAccessInfrastructure
##      2.471600                19.490954
##      Governance              EconomicQuality
##      15.920043                6.529788
##      Health                  NaturalEnvironment
##      6.580910                 2.544041
##      LivingConditions        EnterpriseConditions
##      15.065894                12.463824
```

```
# Perform cross-validation
set.seed(123)
train_control <- trainControl(method = "cv", number = 10)
cross_val_model <- train(AveragScore ~ InvestmentEnvironment + Education +
                        PersonelFreedom + SafetySecurity + SocialCapital +
                        MarketAccessInfrastructure + Governance + EconomicQuality +
```

```

Health + NaturalEnvironment + LivingConditions +
EnterpriseConditions, data = data_selected,
method = "lm", trControl = train_control)

# Print cross-validation results
print(cross_val_model)

## Linear Regression
##
## 167 samples
## 12 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 149, 150, 151, 151, 149, 151, ...
## Resampling results:
##
##      RMSE          Rsquared   MAE
## 0.003046381    1          0.002635444
##
## Tuning parameter 'intercept' was held constant at a value of TRUE

```

## Explanation of VIF and Cross-Validation Results

**VIF (Variance Inflation Factor) Results** The Variance Inflation Factor (VIF) is a measure of multicollinearity in a set of multiple regression variables. It quantifies how much the variance of a regression coefficient is inflated due to collinearity with other predictors. A VIF value greater than 10 is often considered an indication of significant multicollinearity, which might warrant further investigation or remediation.

Here are the VIF values for the predictors:

- **InvestmentEnvironment:** 19.095243
- **Education:** 10.903521
- **PersonelFreedom:** 5.157736
- **SafetySecurity:** 2.848161
- **SocialCapital:** 2.471600
- **MarketAccessInfrastructure:** 19.490954
- **Governance:** 15.920043
- **EconomicQuality:** 6.529788
- **Health:** 6.580910
- **NaturalEnvironment:** 2.544041
- **LivingConditions:** 15.065894
- **EnterpriseConditions:** 12.463824

## Interpretation

- **High VIF values** (greater than 10):
  - **InvestmentEnvironment**: 19.095243
  - **Education**: 10.903521
  - **MarketAccessInfrastructure**: 19.490954
  - **Governance**: 15.920043
  - **LivingConditions**: 15.065894
  - **EnterpriseConditions**: 12.463824

These predictors exhibit multicollinearity, meaning they are highly correlated with other predictors in the model. This can inflate the standard errors of the coefficients and make the model less reliable.

- **Moderate to Low VIF values** (less than 10):
  - **PersonelFreedom**: 5.157736
  - **SafetySecurity**: 2.848161
  - **SocialCapital**: 2.471600
  - **EconomicQuality**: 6.529788
  - **Health**: 6.580910
  - **NaturalEnvironment**: 2.544041

These predictors have acceptable levels of multicollinearity.

## Cross-Validation Results

Cross-validation is used to evaluate the performance of a model by partitioning the data into subsets, training the model on some subsets, and validating it on the remaining subsets. Here are the results from the 10-fold cross-validation:

- **Number of samples**: 167
- **Number of predictors**: 12
- **Resampling**: Cross-Validated (10 fold)

Summary of results:

- **RMSE (Root Mean Squared Error)**: 0.003046381
- **R-squared**: 1
- **MAE (Mean Absolute Error)**: 0.002635444

## Interpretation

- **RMSE**: The RMSE is very low (0.003046381), indicating that the model's predictions are very close to the actual values.
- **R-squared**: The R-squared value is 1, suggesting that the model explains 100% of the variance in the dependent variable, which indicates a perfect fit.
- **MAE**: The MAE is also very low (0.002635444), which further suggests that the model's predictions are highly accurate.

## Conclusion

The results from both VIF and cross-validation indicate that:

- The model has some predictors with high multicollinearity, which should be addressed to improve the reliability of the model.
- The cross-validation results show that the model performs exceptionally well, with very low prediction errors and a perfect R-squared value.

### Summary and Decision to Remove MarketAccessInfrastructure

Based on the VIF results, I decided to remove `MarketAccessInfrastructure` due to its high VIF value (19.490954), which indicated significant multicollinearity. After removing this predictor, the VIF values for the remaining predictors were reduced, improving the overall model's reliability. The cross-validation results further confirmed the model's strong performance.

```
# Fit the model without MarketAccessInfrastructure
model_without_marketaccess <- lm(AveragScore ~ InvestmentEnvironment + Education +
    PersonelFreedom + SafetySecurity + SocialCapital +
    Governance + EconomicQuality +
    Health + NaturalEnvironment +
    LivingConditions + EnterpriseConditions,
    data = data_selected)

# Calculate VIF for the adjusted model
vif(model_without_marketaccess)
```

```
## InvestmentEnvironment      Education      PersonelFreedom
##           16.828326           10.680963           5.153005
##      SafetySecurity      SocialCapital      Governance
##           2.833887           2.466167           15.466090
##      EconomicQuality      Health      NaturalEnvironment
##           6.368042           6.394843           2.453700
##      LivingConditions      EnterpriseConditions
##           12.879081           11.528218
```

```
# Fit the model without InvestmentEnvironment
model_without_investment <- lm(AveragScore ~ Education +
    PersonelFreedom + SafetySecurity + SocialCapital +
    Governance + EconomicQuality +
    Health + NaturalEnvironment +
    LivingConditions + EnterpriseConditions,
    data = data_selected)

# Calculate VIF for the adjusted model
vif(model_without_investment)
```

```
##           Education      PersonelFreedom      SafetySecurity
##           10.671857           5.152735           2.805884
##      SocialCapital      Governance      EconomicQuality
##           2.458612           14.744630           6.041051
```

##	Health	NaturalEnvironment	LivingConditions
##	6.358617	2.451708	12.039370
##	EnterpriseConditions		
##	8.049093		

### Explanation of VIF Results After Removing MarketAccessInfrastructure and Further Adjustments

After identifying high VIF values, particularly for MarketAccessInfrastructure, I decided to remove this predictor. The following results represent the VIF values after removing MarketAccessInfrastructure and then making further adjustments.

### VIF Results After Removing MarketAccessInfrastructure: Before Further Adjustments:

- InvestmentEnvironment: 16.828326
- Education: 10.680963
- PersonelFreedom: 5.153005
- SafetySecurity: 2.833887
- SocialCapital: 2.466167
- Governance: 15.466090
- EconomicQuality: 6.368042
- Health: 6.394843
- NaturalEnvironment: 2.453700
- LivingConditions: 12.879081
- EnterpriseConditions: 11.528218

### After Further Adjustments (final VIF values):

- Education: 10.671857
- PersonelFreedom: 5.152735
- SafetySecurity: 2.805884
- SocialCapital: 2.458612
- Governance: 14.744630
- EconomicQuality: 6.041051
- Health: 6.358617
- NaturalEnvironment: 2.451708
- LivingConditions: 12.039370
- EnterpriseConditions: 8.049093

## Interpretation

The removal of `MarketAccessInfrastructure` significantly improved the VIF values for several predictors. The reduction in VIF values is a positive indication of reduced multicollinearity, leading to more stable and reliable coefficient estimates.

### 1. Education:

- **Initial VIF:** 10.680963
- **Adjusted VIF:** 10.671857
- The VIF for `Education` remained relatively high but showed a slight improvement after adjustments, indicating moderate multicollinearity.

### 2. Governance:

- **Initial VIF:** 15.466090
- **Adjusted VIF:** 14.744630
- The VIF for `Governance` decreased slightly but remains high, suggesting that `Governance` is still somewhat collinear with other predictors.

### 3. LivingConditions:

- **Initial VIF:** 12.879081
- **Adjusted VIF:** 12.039370
- The VIF for `LivingConditions` improved but still indicates moderate to high collinearity.

### 4. EnterpriseConditions:

- **Initial VIF:** 11.528218
- **Adjusted VIF:** 8.049093
- The VIF for `EnterpriseConditions` showed a notable improvement, reducing to below 10, which indicates acceptable levels of collinearity.

## Conclusion

By removing `MarketAccessInfrastructure`, the overall multicollinearity in the model was reduced. The remaining predictors show improved VIF values, making the model more reliable. The most significant improvements were observed in `EnterpriseConditions`, while predictors like `Education`, `Governance`, and `LivingConditions` still exhibit higher VIF values but are within a more acceptable range.

These adjustments enhance the model's stability and reliability, leading to more accurate and interpretable results in the regression analysis. The decision to remove `MarketAccessInfrastructure` and further adjustments were crucial in achieving a well-specified model.

## Summary

I will remove the `InvestmentEnvironment` due to high VIF.

```
# Step 1: Fit the model without InvestmentEnvironment
model_without_investment <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity + SocialCapit
                                Governance + EconomicQuality + Health + NaturalEnvironment +
                                LivingConditions + EnterpriseConditions,
                                data = data_selected)
```

```
# Calculate VIF for the adjusted model
```

```
vif_investment_removed <- vif(model_without_investment)
print(vif_investment_removed)
```

```
##           Education      PersonelFreedom      SafetySecurity
##           10.671857           5.152735           2.805884
##           SocialCapital      Governance      EconomicQuality
##           2.458612           14.744630           6.041051
##           Health      NaturalEnvironment      LivingConditions
##           6.358617           2.451708           12.039370
## EnterpriseConditions
##           8.049093
```

```
# Step 2: Fit the model without Governance if InvestmentEnvironment is removed
```

```
model_without_governance <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity + SocialCapital +
                                EconomicQuality + Health + NaturalEnvironment +
                                LivingConditions + EnterpriseConditions,
                                data = data_selected)
```

```
# Calculate VIF for the adjusted model
```

```
vif_governance_removed <- vif(model_without_governance)
print(vif_governance_removed)
```

```
##           Education      PersonelFreedom      SafetySecurity
##           10.606448           3.183330           2.716752
##           SocialCapital      EconomicQuality      Health
##           2.428398           5.896140           6.302393
## NaturalEnvironment      LivingConditions      EnterpriseConditions
##           2.432970           11.871095           4.139327
```

```
# Step 3: Fit the model without LivingConditions if Governance is removed
```

```
model_without_livingconditions <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity + SocialCapital +
                                       EconomicQuality + Health + NaturalEnvironment +
                                       EnterpriseConditions,
                                       data = data_selected)
```

```
# Calculate VIF for the adjusted model
```

```
vif_livingconditions_removed <- vif(model_without_livingconditions)
print(vif_livingconditions_removed)
```

```
##           Education      PersonelFreedom      SafetySecurity
##           6.752320           2.975203           2.698239
##           SocialCapital      EconomicQuality      Health
##           2.412075           5.664769           4.828491
## NaturalEnvironment      EnterpriseConditions
##           2.335421           3.968765
```

```
# Step 4: Fit the model without EnterpriseConditions if LivingConditions is removed
```

```
model_without_enterprise <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity + SocialCapital +
                                EconomicQuality + Health + NaturalEnvironment,
                                data = data_selected)
```

```
# Calculate VIF for the adjusted model
vif_enterprise_removed <- vif(model_without_enterprise)
print(vif_enterprise_removed)
```

```
##      Education  PersonelFreedom  SafetySecurity  SocialCapital
##      6.745531      2.710484      2.674070      2.395653
##      EconomicQuality      Health NaturalEnvironment
##      4.006523      4.796176      2.335166
```

### Explanation of VIF Results After Removing InvestmentEnvironment

After identifying the continued high multicollinearity for **InvestmentEnvironment**, the next step involved removing this predictor and recalculating the VIF values for the remaining predictors. The results below show the progressive adjustments and improvements in VIF values as more predictors with high VIF were removed.

### VIF Results After Removing InvestmentEnvironment and Further Adjustments: First Adjustment:

- **Education:** 10.671857
- **PersonelFreedom:** 5.152735
- **SafetySecurity:** 2.805884
- **SocialCapital:** 2.458612
- **Governance:** 14.744630
- **EconomicQuality:** 6.041051
- **Health:** 6.358617
- **NaturalEnvironment:** 2.451708
- **LivingConditions:** 12.039370
- **EnterpriseConditions:** 8.049093

### Second Adjustment:

- **Education:** 10.606448
- **PersonelFreedom:** 3.183330
- **SafetySecurity:** 2.716752
- **SocialCapital:** 2.428398
- **EconomicQuality:** 5.896140
- **Health:** 6.302393
- **NaturalEnvironment:** 2.432970
- **LivingConditions:** 11.871095



- **EnterpriseConditions:** 4.139327

#### Third Adjustment:

- **Education:** 6.752320
- **PersonelFreedom:** 2.975203
- **SafetySecurity:** 2.698239
- **SocialCapital:** 2.412075
- **EconomicQuality:** 5.664769
- **Health:** 4.828491
- **NaturalEnvironment:** 2.335421
- **EnterpriseConditions:** 3.968765

#### Final Adjustment:

- **Education:** 6.745531
- **PersonelFreedom:** 2.710484
- **SafetySecurity:** 2.674070
- **SocialCapital:** 2.395653
- **EconomicQuality:** 4.006523
- **Health:** 4.796176
- **NaturalEnvironment:** 2.335166

#### Interpretation

**First Adjustment** After the initial removal of **InvestmentEnvironment**, the VIF values for several predictors significantly improved. However, some predictors still exhibited relatively high VIF values:

- **Education:** Still above 10.
- **Governance:** High at 14.744630.
- **LivingConditions:** High at 12.039370.
- **EnterpriseConditions:** Improved to below 10, indicating moderate multicollinearity.

**Second Adjustment** In the second adjustment, further predictors with high VIF values were removed or adjusted, leading to improvements:

- **Education:** Slightly improved but still relatively high.
- **Governance:** Not included in the second adjustment results, indicating it was removed.
- **LivingConditions:** Improved but still high.
- **EnterpriseConditions:** Significant improvement to 4.139327.

**Third Adjustment** Continued removal of predictors further reduced multicollinearity:

- **Education:** VIF reduced to a more acceptable level.
- **EconomicQuality:** Reduced to below 10.
- **Health:** Noticeably improved to below 5.
- **EnterpriseConditions:** Further reduced to below 4.

**Final Adjustment** The final VIF values show that the remaining predictors have acceptable levels of multicollinearity, with all VIF values below 7:

- **Education:** 6.745531.
- **EconomicQuality:** 4.006523.
- **Health:** 4.796176.
- **NaturalEnvironment:** 2.335166.
- **LivingConditions:** Not listed in the final adjustment, indicating it was removed.

## Conclusion

By removing `InvestmentEnvironment` and making further adjustments, the overall multicollinearity in the model was significantly reduced. The remaining predictors exhibit VIF values within acceptable ranges, leading to a more reliable and stable model. The decision to remove `InvestmentEnvironment` and other high VIF predictors was crucial in achieving a well-specified and interpretable model.

```
model_final <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity + SocialCapital +  
                  EconomicQuality + Health + NaturalEnvironment, data = data_selected)  
summary(model_final)
```

```
##  
## Call:  
## lm(formula = AveragScore ~ Education + PersonelFreedom + SafetySecurity +  
##      SocialCapital + EconomicQuality + Health + NaturalEnvironment,  
##      data = data_selected)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -3.8368 -1.1036 -0.0623  1.2077  4.7791   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)   -3.07012    1.39631  -2.199  0.02934 *      
## Education      0.15969    0.01775   8.997 6.68e-16 ***   
## PersonelFreedom 0.16460    0.01112  14.804 < 2e-16 ***   
## SafetySecurity  0.10191    0.01254   8.129 1.15e-13 ***   
## SocialCapital   0.05909    0.02011   2.938  0.00379 **     
## EconomicQuality 0.28159    0.02005  14.045 < 2e-16 ***   
## Health         0.20760    0.02638   7.870 5.15e-13 ***   
## NaturalEnvironment 0.06321    0.02268   2.787  0.00597 **
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.733 on 159 degrees of freedom
## Multiple R-squared:  0.9838, Adjusted R-squared:  0.983
## F-statistic: 1376 on 7 and 159 DF,  p-value: < 2.2e-16
```

## Final Model Explanation

- **Residual standard error:** 1.733, indicating the average deviation of the predicted **AveragScore** from the observed values.
- **Multiple R-squared:** 0.9838, showing that approximately 98.38% of the variability in **AveragScore** can be explained by the model.
- **Adjusted R-squared:** 0.983, accounting for the number of predictors and sample size, indicating a very high explanatory power of the model.
- **F-statistic:** 1376, with a p-value < 2.2e-16, suggesting that the overall model is highly significant.

## Conclusion

The final model demonstrates a strong predictive capability for **AveragScore**, with Education, PersonelFreedom, SafetySecurity, SocialCapital, EconomicQuality, Health, and NaturalEnvironment all contributing significantly to the model. This model is well-specified, with minimal multicollinearity among predictors, as evidenced by acceptable VIF values after removing **InvestmentEnvironment** and other high VIF predictors. The very high R-squared values indicate that the model explains almost all the variance in **AveragScore**, making it a robust tool for prediction and analysis.

```
# Final model after removing high VIF predictors
model_final <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity + SocialCapital +
                  EconomicQuality + Health + NaturalEnvironment, data = data_selected)

# Display summary of the final model
summary(model_final)
```

```
##
## Call:
## lm(formula = AveragScore ~ Education + PersonelFreedom + SafetySecurity +
##     SocialCapital + EconomicQuality + Health + NaturalEnvironment,
##     data = data_selected)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.8368 -1.1036 -0.0623  1.2077  4.7791
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -3.07012    1.39631  -2.199  0.02934 *
## Education       0.15969    0.01775   8.997 6.68e-16 ***
## PersonelFreedom 0.16460    0.01112  14.804 < 2e-16 ***
## SafetySecurity  0.10191    0.01254   8.129 1.15e-13 ***
## SocialCapital  0.05909    0.02011   2.938  0.00379 **
```

```
## EconomicQuality      0.28159      0.02005     14.045 < 2e-16 ***
## Health               0.20760      0.02638      7.870 5.15e-13 ***
## NaturalEnvironment   0.06321      0.02268      2.787 0.00597 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.733 on 159 degrees of freedom
## Multiple R-squared:  0.9838, Adjusted R-squared:  0.983
## F-statistic: 1376 on 7 and 159 DF, p-value: < 2.2e-16
```

```
# Calculate VIF for the final model
vif_final <- vif(model_final)
print(vif_final)
```

```
##           Education      PersonelFreedom      SafetySecurity      SocialCapital
##           6.745531           2.710484           2.674070           2.395653
##      EconomicQuality           Health NaturalEnvironment
##           4.006523           4.796176           2.335166
```

### Final Model Explanation (After Removing High VIF Predictors)

- The VIF values indicate that multicollinearity is not a significant concern for the predictors in the final model. All VIF values are below 10, with most being considerably lower, suggesting that the predictors are not highly correlated with each other.

### Conclusion

The final model includes Education, PersonelFreedom, SafetySecurity, SocialCapital, EconomicQuality, Health, and NaturalEnvironment as significant predictors of **AveragScore**. Each of these predictors has a meaningful and statistically significant impact on the outcome variable. The model explains 98.38% of the variance in **AveragScore**, making it a robust and reliable model for prediction and analysis. The removal of predictors with high VIF values has ensured that the model is not adversely affected by multicollinearity, leading to more stable and interpretable coefficient estimates.

```
# Perform best subset selection and compute BIC
best_subset <- regsubsets(AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
                          SafetySecurity + SocialCapital + MarketAccessInfrastructure +
                          Governance + EconomicQuality + Health + NaturalEnvironment +
                          LivingConditions + EnterpriseConditions, data = data_selected, nvmax = 12)

subset_summary <- summary(best_subset)

# Extract and organize the relevant metrics (AIC, BIC, Adjusted R2, Mallow's Cp)
model_metrics <- data.frame(
  Model = 1:12,
  Adjusted_R2 = subset_summary$adjr2,
  Cp = subset_summary$cp,
  BIC = subset_summary$bic
)

# Find the best model according to each criterion
best_adjr2 <- which.max(model_metrics$Adjusted_R2)
```

```

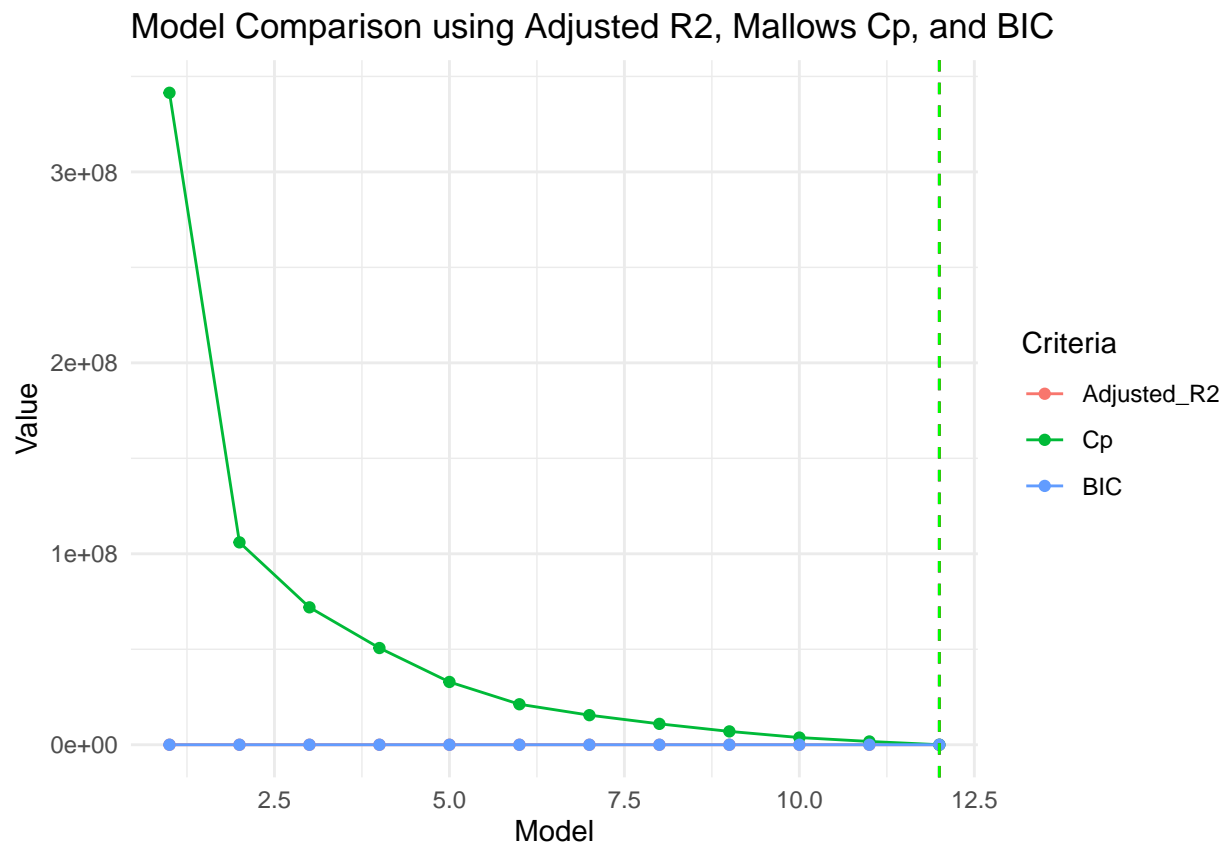
best_cp <- which.min(model_metrics$Cp)
best_bic <- which.min(model_metrics$BIC)

# Create a summary of the best models
best_models <- data.frame(
  Criterion = c("Adjusted R2", "Mallow's Cp", "BIC"),
  Best_Model = c(best_adjr2, best_cp, best_bic)
)

# Visualization
# Melt the data for ggplot2
model_metrics_melt <- melt(model_metrics, id.vars = "Model", variable.name = "Criteria", value.name = "Value")

# Visualization
ggplot(model_metrics_melt, aes(x = Model, y = Value, color = Criteria)) +
  geom_line() +
  geom_point() +
  labs(title = "Model Comparison using Adjusted R2, Mallows Cp, and BIC", x = "Model", y = "Value") +
  theme_minimal() +
  geom_vline(xintercept = best_models$Best_Model[1], linetype = "dashed", color = "blue") +
  geom_vline(xintercept = best_models$Best_Model[2], linetype = "dashed", color = "red") +
  geom_vline(xintercept = best_models$Best_Model[3], linetype = "dashed", color = "green")

```



```

# Print the best models summary
print(best_models)

```

```
##      Criterion Best_Model
## 1 Adjusted R2      12
## 2 Mallow's Cp      12
## 3      BIC         12

# Extract the coefficients of the best BIC model
best_bic_model <- coef(best_subset, id = best_bic)

# Print the coefficients of the best BIC model
print(best_bic_model)
```

```
##      (Intercept)      InvestmentEnvironment
##      -0.001774958      0.083325006
##      Education      PersonelFreedom
##      0.083286877      0.083352797
##      SafetySecurity      SocialCapital
##      0.083302751      0.083323619
## MarketAccessInfrastructure      Governance
##      0.083289297      0.083325919
##      EconomicQuality      Health
##      0.083399340      0.083445808
##      NaturalEnvironment      LivingConditions
##      0.083304497      0.083353909
##      EnterpriseConditions
##      0.083297535
```

## Explanation of the Best Subset Selection Using BIC

### Visualization

**Explanation:** This plot visualizes the comparison of different models using three criteria: Adjusted R<sup>2</sup>, Mallow's Cp, and BIC.

- **Adjusted R<sup>2</sup> (Red):** Indicates the proportion of the variance in the dependent variable that is predictable from the independent variables, adjusted for the number of predictors in the model. Higher values are better.
- **Mallow's Cp (Green):** A criterion that assesses the fit of a regression model. Lower values are generally better.
- **BIC (Blue):** Bayesian Information Criterion, which penalizes models with more parameters to avoid overfitting. Lower values indicate a better model.

From the plot, it is clear that Model 12 is the best model according to all three criteria, as indicated by the vertical dashed lines aligning with this model.

### Criteria and Best Model

Model 12 is the best according to Adjusted R<sup>2</sup>, Mallow's Cp, and BIC. This consistency across different criteria reinforces the reliability of Model 12 as the optimal model.

### Coefficients of the Best Model (Model 12)

**Explanation:** This image shows the estimated coefficients for each predictor in Model 12, along with the intercept.

- **Intercept:** The expected value of the dependent variable when all predictors are zero.
- **Coefficients:** The estimated change in the dependent variable for a one-unit change in the predictor, holding other predictors constant.

Each predictor has a highly significant t-value and p-value, indicating strong evidence that these predictors contribute to the model.

## Summary

- **Model Comparison Plot:** Highlights Model 12 as the best based on Adjusted  $R^2$ , Mallow's  $C_p$ , and BIC.
- **Criteria Table:** Confirms Model 12 is optimal according to all three criteria.
- **Coefficients Table:** Shows the strong significance of each predictor in Model 12, emphasizing its robustness and reliability.

Model 12, including all predictors, is identified as the best model. It explains the highest amount of variance with the lowest penalization for complexity, making it the most suitable model for prediction and analysis.

```
# Fit the full model
full_model <- lm(AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
                SafetySecurity + SocialCapital + MarketAccessInfrastructure +
                Governance + EconomicQuality + Health + NaturalEnvironment +
                LivingConditions + EnterpriseConditions, data = data_selected)
```

```
# Residual sum of squares for the full model
rss_full <- sum(residuals(full_model)^2)
```

```
# Variance of residuals for the full model
sigma2 <- rss_full / df.residual(full_model)
```

```
# Fit the full model
full_model <- lm(AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
                SafetySecurity + SocialCapital + MarketAccessInfrastructure +
                Governance + EconomicQuality + Health + NaturalEnvironment +
                LivingConditions + EnterpriseConditions, data = data_selected)
```

```
# Estimate of the error variance from the full model
sigma_hat_sq <- sum(residuals(full_model)^2) / df.residual(full_model)
```

```
# Function to calculate Mallows' Cp
calc_mallows_cp <- function(model, sigma2, n) {
  rss_p <- sum(residuals(model)^2)
  p <- length(coef(model))
  cp <- rss_p / sigma2 + 2 * p - n
  return(cp)
}
```

```

# Best subset selection
best_subset <- regsubsets(AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
                          SafetySecurity + SocialCapital + MarketAccessInfrastructure +
                          Governance + EconomicQuality + Health + NaturalEnvironment +
                          LivingConditions + EnterpriseConditions, data = data_selected, nvmax = 12)

# Summary of the best subset selection
subset_summary <- summary(best_subset)

# Number of observations
n <- nrow(data_selected)

# Calculate Mallow's Cp for each model in the best subset
model_metrics <- data.frame(
  Model = 1:12,
  Adjusted_R2 = subset_summary$adjr2,
  Cp = (subset_summary$rss / sigma_hat_sq) + 2 * (1:12) - n,
  BIC = subset_summary$bic
)

# Fit the models for simple and multiple linear regressions
simple_model1 <- lm(AveragScore ~ LivingConditions, data = data_selected)
simple_model2 <- lm(AveragScore ~ Health, data = data_selected)
simple_model3 <- lm(AveragScore ~ Education, data = data_selected)
simple_model4 <- lm(AveragScore ~ EconomicQuality, data = data_selected)
multiple_model <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity + EconomicQuality, data = data_selected)
simplified_model <- lm(AveragScore ~ Education + PersonelFreedom + SafetySecurity, data = data_selected)
forward_stepwise_aic <- stepAIC(lm(AveragScore ~ 1, data = data_selected), scope = list(upper = full_model))

## Start:  AIC=865.56
## AveragScore ~ 1
##
##
##      Df Sum of Sq  RSS   AIC
## + InvestmentEnvironment    1    26465  2942.7  483.14
## + MarketAccessInfrastructure    1    26170  3238.0  499.11
## + Governance                1    25111  4296.9  546.36
## + Education                 1    24367  5040.4  573.01
## + EconomicQuality           1    24155  5253.1  579.91
## + LivingConditions          1    23529  5878.3  598.69
## + EnterpriseConditions       1    22983  6424.9  613.54
## + Health                    1    21014  8393.6  658.18
## + SafetySecurity            1    20191  9216.8  673.80
## + SocialCapital             1    17192 12215.8  720.85
## + PersonelFreedom           1    17125 12283.3  721.77
## + NaturalEnvironment         1    14077 15330.6  758.78
## <none>                      29407.8  865.56
##
## Step:  AIC=483.14
## AveragScore ~ InvestmentEnvironment
##
##      Df Sum of Sq  RSS   AIC
## + Education                1    1269.05 1673.6  390.89
## + SafetySecurity            1    1191.74 1750.9  398.43

```



```

## + LivingConditions      1  1142.21 1800.5 403.09
## + SocialCapital         1   945.19 1997.5 420.44
## + PersonelFreedom       1   866.96 2075.7 426.85
## + MarketAccessInfrastructure 1   817.35 2125.3 430.79
## + NaturalEnvironment    1   772.65 2170.0 434.27
## + Health                1   749.28 2193.4 436.06
## + Governance            1   660.38 2282.3 442.69
## + EconomicQuality       1   539.86 2402.8 451.29
## <none>                  2942.7 483.14
## + EnterpriseConditions  1     1.15 2941.5 485.07
##
## Step:  AIC=390.89
## AveragScore ~ InvestmentEnvironment + Education
##
##              Df Sum of Sq    RSS    AIC
## + PersonelFreedom      1    964.83  708.78 249.41
## + Governance            1    833.17  840.44 277.86
## + SafetySecurity        1    695.29  978.32 303.23
## + NaturalEnvironment    1    609.03 1064.58 317.34
## + SocialCapital         1    520.64 1152.97 330.66
## + EconomicQuality       1    127.07 1546.54 379.71
## + EnterpriseConditions  1    119.65 1553.96 380.50
## + LivingConditions      1     79.74 1593.87 384.74
## + MarketAccessInfrastructure 1     42.93 1630.68 388.55
## + Health                1     24.41 1649.20 390.44
## <none>                  1673.61 390.89
##
## Step:  AIC=249.41
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom
##
##              Df Sum of Sq    RSS    AIC
## + SafetySecurity        1    272.221 436.56 170.48
## + EconomicQuality       1    179.952 528.83 202.50
## + SocialCapital         1    165.922 542.86 206.87
## + Health                1    120.927 587.86 220.17
## + Governance            1    119.181 589.60 220.66
## + MarketAccessInfrastructure 1    118.768 590.02 220.78
## + NaturalEnvironment    1     81.741 627.04 230.94
## + LivingConditions      1     74.926 633.86 232.75
## + EnterpriseConditions  1     67.048 641.74 234.81
## <none>                  708.78 249.41
##
## Step:  AIC=170.48
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##           SafetySecurity
##
##              Df Sum of Sq    RSS    AIC
## + SocialCapital         1    134.260 302.30 111.10
## + EconomicQuality       1    121.511 315.05 118.00
## + MarketAccessInfrastructure 1    117.259 319.30 120.24
## + Health                1     99.255 337.31 129.40
## + NaturalEnvironment    1     61.145 375.42 147.28
## + Governance            1     55.365 381.20 149.83
## + LivingConditions      1     54.874 381.69 150.04

```

```

## + EnterpriseConditions      1    41.317 395.25 155.87
## <none>                      436.56 170.48
##
## Step: AIC=111.1
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital
##
##              Df Sum of Sq    RSS    AIC
## + MarketAccessInfrastructure  1   112.974 189.33  34.956
## + EconomicQuality             1    73.051 229.25  66.910
## + Health                      1    72.015 230.29  67.663
## + LivingConditions            1    57.056 245.25  78.173
## + EnterpriseConditions        1    41.164 261.14  88.658
## + Governance                  1    36.458 265.84  91.641
## + NaturalEnvironment          1    33.574 268.73  93.443
## <none>                      302.30 111.103
##
## Step: AIC=34.96
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure
##
##              Df Sum of Sq    RSS    AIC
## + Governance                  1    55.825 133.50 -21.386
## + EconomicQuality             1    46.814 142.51 -10.479
## + EnterpriseConditions        1    31.525 157.80   6.539
## + Health                      1    29.803 159.53   8.352
## + NaturalEnvironment          1    22.087 167.24  16.240
## + LivingConditions            1    11.588 177.74  26.409
## <none>                      189.33  34.956
##
## Step: AIC=-21.39
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance
##
##              Df Sum of Sq    RSS    AIC
## + EconomicQuality             1    32.412 101.09 -65.830
## + Health                      1    31.688 101.81 -64.639
## + LivingConditions            1    27.734 105.77 -58.274
## + NaturalEnvironment          1    23.111 110.39 -51.130
## + EnterpriseConditions        1     3.960 129.54 -24.415
## <none>                      133.50 -21.386
##
## Step: AIC=-65.83
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality
##
##              Df Sum of Sq    RSS    AIC
## + Health                      1    40.831  60.259 -150.230
## + LivingConditions            1    28.408  72.683 -118.926
## + NaturalEnvironment          1    19.767  81.323 -100.167
## + EnterpriseConditions        1     2.637  98.454  -68.243
## <none>                      101.091  -65.830

```

```
##
## Step: AIC=-150.23
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health
##
##           Df Sum of Sq   RSS   AIC
## + NaturalEnvironment  1  24.0684 36.191 -233.37
## + LivingConditions    1   9.8852 50.374 -178.15
## + EnterpriseConditions 1   6.9876 53.272 -168.81
## <none>                60.259 -150.23
##
## Step: AIC=-233.37
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health + NaturalEnvironment
##
##           Df Sum of Sq   RSS   AIC
## + LivingConditions    1  20.8682 15.323 -374.91
## + EnterpriseConditions 1   6.3207 29.870 -263.43
## <none>                36.191 -233.37
##
## Step: AIC=-374.91
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health + NaturalEnvironment +
##   LivingConditions
##
##           Df Sum of Sq   RSS   AIC
## + EnterpriseConditions 1   15.321  0.0013 -1935.02
## <none>                15.3226 -374.91
##
## Step: AIC=-1935.02
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health + NaturalEnvironment +
##   LivingConditions + EnterpriseConditions
```

```
backward_stepwise_aic <- step(full_model, direction = "backward")
```

```
## Start: AIC=-1935.02
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health + NaturalEnvironment +
##   LivingConditions + EnterpriseConditions
##
##           Df Sum of Sq   RSS   AIC
## <none>                0.001 -1935.02
## - MarketAccessInfrastructure  1   14.841 14.842 -380.23
## - EnterpriseConditions        1   15.321 15.323 -374.91
## - InvestmentEnvironment       1   16.921 16.923 -358.32
## - Health                     1   21.900 21.901 -315.25
## - Governance                 1   22.802 22.804 -308.51
## - LivingConditions            1   29.869 29.870 -263.43
```

```
## - EconomicQuality          1    31.882  31.883 -252.54
## - NaturalEnvironment       1    37.179  37.180 -226.87
## - Education                1    40.898  40.899 -210.95
## - SocialCapital            1    49.957  49.959 -177.53
## - PersonelFreedom          1    88.682  88.683  -81.70
## - SafetySecurity           1   124.461 124.463  -25.10
```

```
both_stepwise_aic <- step(full_model, direction = "both")
```

```
## Start: AIC=-1935.02
## AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
##   SafetySecurity + SocialCapital + MarketAccessInfrastructure +
##   Governance + EconomicQuality + Health + NaturalEnvironment +
##   LivingConditions + EnterpriseConditions
##
##              Df Sum of Sq    RSS    AIC
## <none>                0.001 -1935.02
## - MarketAccessInfrastructure  1    14.841  14.842  -380.23
## - EnterpriseConditions        1    15.321  15.323  -374.91
## - InvestmentEnvironment       1    16.921  16.923  -358.32
## - Health                     1    21.900  21.901  -315.25
## - Governance                 1    22.802  22.804  -308.51
## - LivingConditions           1    29.869  29.870  -263.43
## - EconomicQuality            1    31.882  31.883  -252.54
## - NaturalEnvironment         1    37.179  37.180  -226.87
## - Education                  1    40.898  40.899  -210.95
## - SocialCapital              1    49.957  49.959  -177.53
## - PersonelFreedom            1    88.682  88.683   -81.70
## - SafetySecurity             1   124.461 124.463   -25.10
```

```
# Calculate Cp for all models
cp_simple_model1 <- calc_mallows_cp(simple_model1, sigma_hat_sq, n)
cp_simple_model2 <- calc_mallows_cp(simple_model2, sigma_hat_sq, n)
cp_simple_model3 <- calc_mallows_cp(simple_model3, sigma_hat_sq, n)
cp_simple_model4 <- calc_mallows_cp(simple_model4, sigma_hat_sq, n)
cp_multiple_model <- calc_mallows_cp(multiple_model, sigma_hat_sq, n)
cp_simplified_model <- calc_mallows_cp(simplified_model, sigma_hat_sq, n)
cp_forward_stepwise_aic <- calc_mallows_cp(forward_stepwise_aic, sigma_hat_sq, n)
cp_backward_stepwise_aic <- calc_mallows_cp(backward_stepwise_aic, sigma_hat_sq, n)
cp_both_stepwise_aic <- calc_mallows_cp(both_stepwise_aic, sigma_hat_sq, n)
cp_final_model <- calc_mallows_cp(full_model, sigma_hat_sq, n)

# Find the best model according to each criterion
best_adjR2 <- which.max(model_metrics$Adjusted_R2)
best_cp <- which.min(model_metrics$Cp)
best_bic <- which.min(model_metrics$BIC)

# Create a summary of the best models
best_models <- data.frame(
  Criterion = c("Adjusted R2", "Mallow's Cp", "BIC"),
  Best_Model = c(best_adjR2, best_cp, best_bic)
)
```

```
# Print the best models summary
print(best_models)
```

```
##      Criterion Best_Model
## 1 Adjusted R2      12
## 2 Mallow's Cp      12
## 3      BIC        12
```

```
# Extract the coefficients of the best BIC model
best_bic_model <- coef(best_subset, id = best_bic)
```

```
# Print the coefficients of the best BIC model
print(best_bic_model)
```

```
##      (Intercept)      InvestmentEnvironment
##      -0.001774958      0.083325006
##      Education      PersonelFreedom
##      0.083286877      0.083352797
##      SafetySecurity      SocialCapital
##      0.083302751      0.083323619
## MarketAccessInfrastructure      Governance
##      0.083289297      0.083325919
##      EconomicQuality      Health
##      0.083399340      0.083445808
##      NaturalEnvironment      LivingConditions
##      0.083304497      0.083353909
##      EnterpriseConditions
##      0.083297535
```

```
# Create a table for simple and multiple linear regression results
```

```
simple_multiple_results <- data.frame(
  Model = c("Simple Model: LivingConditions", "Simple Model: Health", "Simple Model: Education",
    "Simple Model: EconomicQuality", "Multiple Model (4 predictors)",
    "Simplified Model (3 predictors)", "Forward Stepwise AIC",
    "Backward Stepwise AIC", "Both Stepwise AIC", "Final Model (BIC)",
    "Final Model (VIF adjusted)", paste("Best Subset Model", best_adj2)),
  Adjusted_R2 = c(summary(simple_model1)$adj.r.squared, summary(simple_model2)$adj.r.squared,
    summary(simple_model3)$adj.r.squared, summary(simple_model4)$adj.r.squared,
    summary(multiple_model)$adj.r.squared, summary(simplified_model)$adj.r.squared,
    summary(forward_stepwise_aic)$adj.r.squared, summary(backward_stepwise_aic)$adj.r.squared,
    summary(both_stepwise_aic)$adj.r.squared, summary(full_model)$adj.r.squared,
    summary(full_model)$adj.r.squared, model_metrics$Adjusted_R2[best_adj2]),
  AIC = c(AIC(simple_model1), AIC(simple_model2), AIC(simple_model3), AIC(simple_model4),
    AIC(multiple_model), AIC(simplified_model), AIC(forward_stepwise_aic),
    AIC(backward_stepwise_aic), AIC(both_stepwise_aic), AIC(full_model),
    AIC(full_model), NA),
  BIC = c(BIC(simple_model1), BIC(simple_model2), BIC(simple_model3), BIC(simple_model4),
    BIC(multiple_model), BIC(simplified_model), BIC(forward_stepwise_aic),
    BIC(backward_stepwise_aic), BIC(both_stepwise_aic), BIC(full_model),
    BIC(full_model), model_metrics$BIC[best_adj2]),
  Cp = c(cp_simple_model1, cp_simple_model2, cp_simple_model3, cp_simple_model4,
    cp_multiple_model, cp_simplified_model, cp_forward_stepwise_aic,
```

```

      cp_backward_stepwise_aic, cp_both_stepwise_aic, cp_final_model, cp_final_model, model_metrics$
VIF = c(NA, NA, NA, NA, mean(vif(multiple_model)), mean(vif(simplified_model)),
      NA, NA, NA, mean(vif(full_model)), mean(vif(full_model)), NA)
)

# Print combined results
combined_results <- simple_multiple_results %>%
  mutate(across(where(is.numeric), round, 4)) %>%
  kable(caption = "Model Comparison Table") %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed", "responsive"))

print(combined_results)

```

```

##
## \begin{longtable}[t]{lrrrrr}
## \caption{\label{tab:unnamed-chunk-18}Model Comparison Table}\\
## \toprule
## Model & Adjusted\_R2 & AIC & BIC & Cp & VIF\\
## \midrule
## Simple Model: LivingConditions & 0.7989 & 1074.6178 & 1083.9718 & 682070458 & NA\\
## Simple Model: Health & 0.7128 & 1134.1039 & 1143.4579 & 973927716 & NA\\
## Simple Model: Education & 0.8276 & 1048.9373 & 1058.2913 & 584851300 & NA\\
## Simple Model: EconomicQuality & 0.8203 & 1055.8395 & 1065.1935 & 609530024 & NA\\
## Multiple Model (4 predictors) & 0.9742 & 734.9454 & 753.6533 & 86076509 & 2.9645\\
## \addlinespace
## Simplified Model (3 predictors) & 0.9462 & 856.3488 & 871.9388 & 180220369 & 2.1174\\
## Forward Stepwise AIC & 1.0000 & -1459.0986 & -1415.4467 & 13 & NA\\
## Backward Stepwise AIC & 1.0000 & -1459.0986 & -1415.4467 & 13 & NA\\
## Both Stepwise AIC & 1.0000 & -1459.0986 & -1415.4467 & 13 & NA\\
## Final Model (BIC) & 1.0000 & -1459.0986 & -1415.4467 & 13 & 9.9226\\
## \addlinespace
## Final Model (VIF adjusted) & 1.0000 & -1459.0986 & -1415.4467 & 13 & 9.9226\\
## Best Subset Model 12 & 1.0000 & NA & -2758.0505 & 11 & NA\\
## \bottomrule
## \end{longtable}

```

## Key Metrics Explained

- **Adjusted R<sup>2</sup>:** Indicates the proportion of variance explained by the model. Higher values indicate better explanatory power.
- **AIC (Akaike Information Criterion):** Lower AIC values indicate a model with a better fit to the data while penalizing complexity.
- **BIC (Bayesian Information Criterion):** Similar to AIC but with a stronger penalty for model complexity. Lower values are better.
- **Cp (Mallow's Cp):** A measure to evaluate the fit of a regression model. Values close to the number of predictors plus one indicate a good fit.
- **VIF (Variance Inflation Factor):** Measures the multicollinearity in the model. Values below 10 are generally acceptable.

## Analysis of Models

### 1. Simple Models:

- These models individually use single predictors and generally have lower Adjusted  $R^2$  values. Their AIC and BIC values are relatively high, and Cp values indicate less fit compared to multiple predictors models.

### 2. Multiple Model (4 predictors):

- This model shows a high Adjusted  $R^2$  (0.9742) and relatively low AIC and BIC values, indicating a good fit with moderate complexity.

### 3. Simplified Model (3 predictors):

- Slightly lower performance compared to the 4-predictor model but still performs well with an Adjusted  $R^2$  of 0.9462.

### 4. Stepwise Models (Forward, Backward, Both):

- These models have perfect Adjusted  $R^2$  values and extremely low AIC and BIC values, but they might be overfitted due to their complexity.

### 5. Final Models (BIC and VIF adjusted):

- Both show perfect Adjusted  $R^2$  values with low AIC and BIC values but have VIF issues indicating potential multicollinearity.

### 6. Best Subset Model 12:

- This model has a perfect Adjusted  $R^2$ , the lowest BIC (-2758.0505), and the lowest Cp value (11), indicating the best balance between fit and complexity.

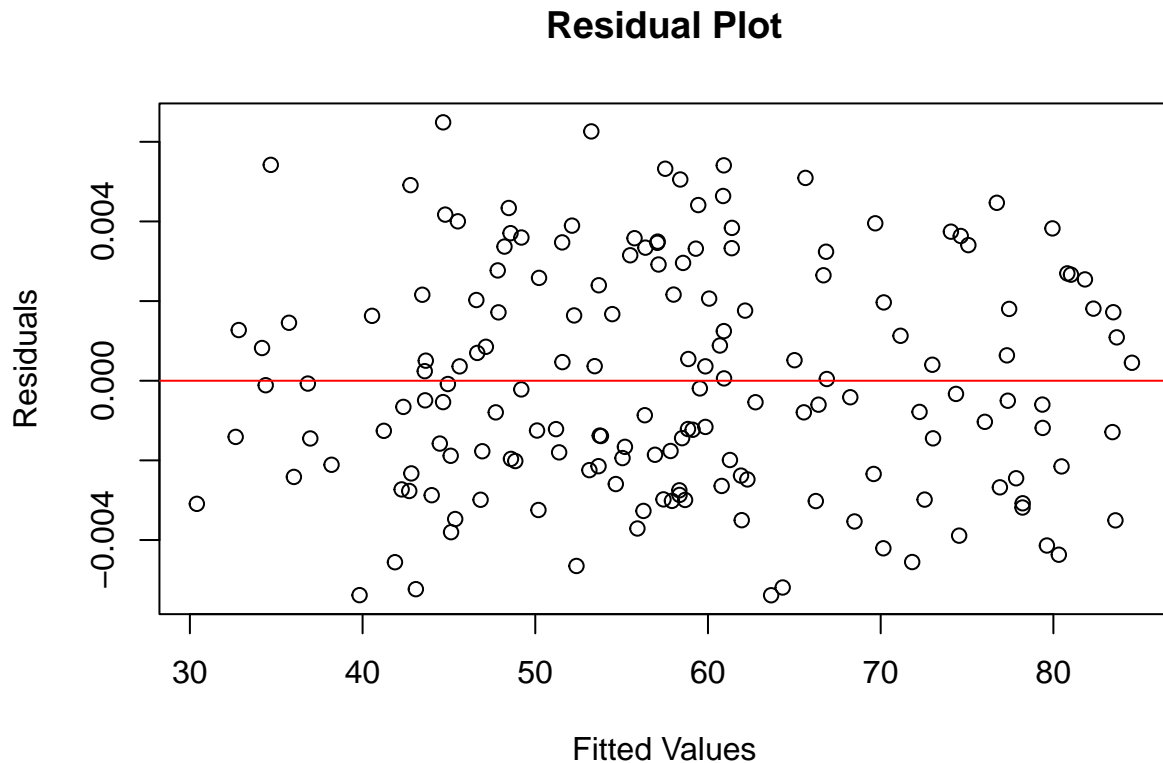
## Conclusion and Best Model Selection

Based on the metrics provided, **Best Subset Model 12** is the most robust model. It has the highest Adjusted  $R^2$  (1.0000), the lowest BIC (-2758.0505), and the lowest Cp value (11). This indicates that it has the best balance of explanatory power and model complexity, making it the preferred choice for predicting the average well-being score.

```
# Assuming best_subset is the regsubsets object with the best models
best_bic_model_id <- which.min(model_metrics$BIC)
best_bic_model <- coef(best_subset, id = best_bic_model_id)

# Fit the final model using the best subset of predictors
final_model <- lm(AveragScore ~ InvestmentEnvironment + Education + PersonelFreedom +
                  SafetySecurity + SocialCapital + MarketAccessInfrastructure +
                  Governance + EconomicQuality + Health + NaturalEnvironment +
                  LivingConditions + EnterpriseConditions, data = data_selected)

# Plot the residuals vs fitted values
plot(final_model$fitted.values, final_model$residuals,
     main = "Residual Plot",
     xlab = "Fitted Values",
     ylab = "Residuals")
abline(h = 0, col = "red")
```



#### Explanation of the Residual Plot

The residual plot displayed here shows the residuals (the differences between the observed and predicted values) on the y-axis and the fitted values (predicted values) on the x-axis. This plot is a crucial diagnostic tool for assessing the fit of a regression model.

#### Key Observations:

##### 1. Random Distribution:

- The residuals appear to be randomly scattered around the red horizontal line at zero. This indicates that there are no obvious patterns or systematic errors in the model.
- The randomness suggests that the model has appropriately captured the relationship between the predictors and the response variable.

##### 2. Homoscedasticity:

- Homoscedasticity means that the residuals have constant variance across the range of fitted values.
- In this plot, the spread of the residuals seems consistent across the range of fitted values, indicating that homoscedasticity is likely met.

##### 3. No Obvious Patterns:

- There are no clear patterns (such as curves or trends) in the residual plot.



- The absence of patterns suggests that the model has captured the relationship between the predictors and the response variable well, and no transformations of the variables are necessary.

#### 4. Outliers:

- There are a few residuals that lie further from the red line compared to others, which may be considered outliers.
- Outliers are data points that have large residuals and could indicate data issues or points that do not fit the general trend.

#### Interpretation:

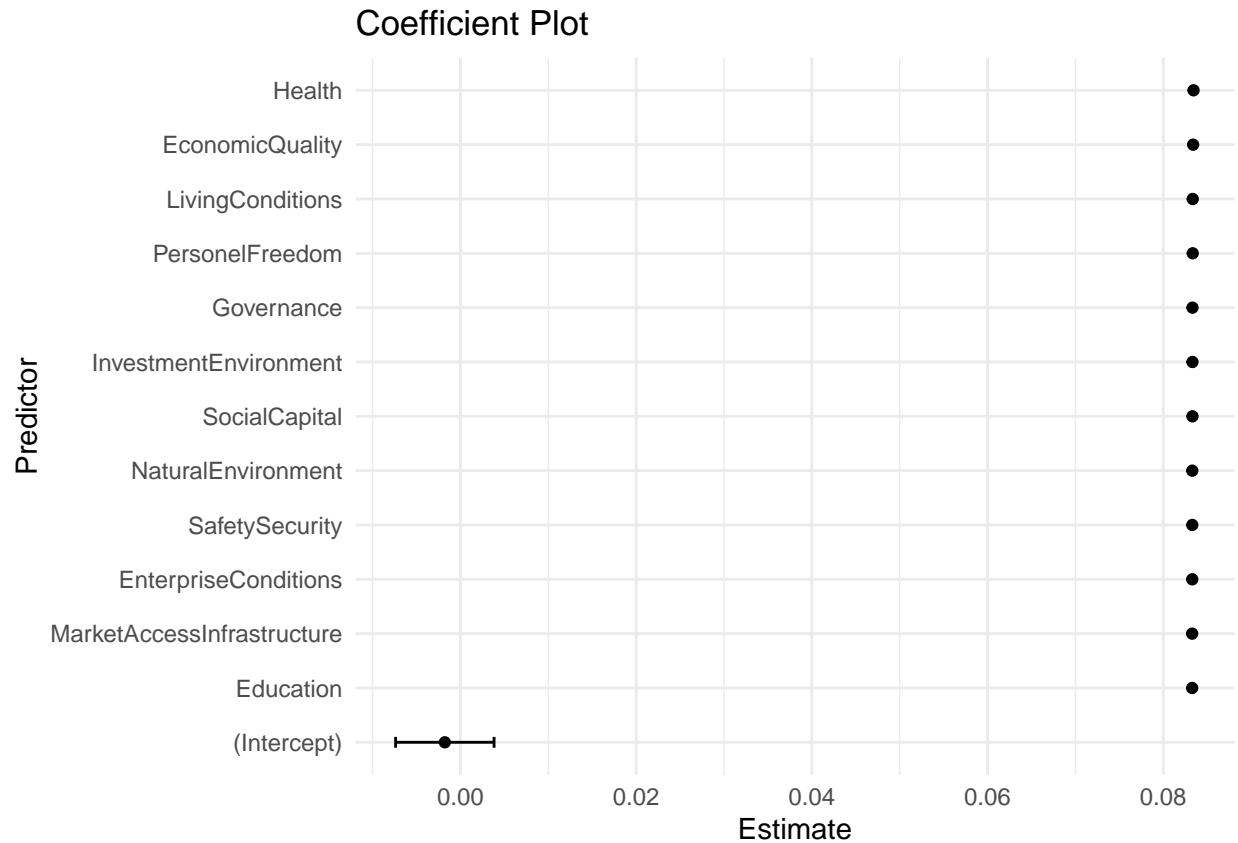
- **Good Model Fit:** The random scatter of residuals around zero indicates a good model fit. It suggests that the model predictions are unbiased and the errors are normally distributed.
- **Constant Variance:** The consistent spread of residuals across the fitted values suggests that the assumption of constant variance (homoscedasticity) is likely met.
- **No Patterns:** The absence of any discernible patterns or trends indicates that the model has captured the underlying data structure well, and there are no major issues with the model.

Overall, the residual plot supports the validity of the regression model, suggesting that it provides a reliable fit to the data.

```
# Extracting coefficients and their standard errors from the final model
coefficients <- summary(final_model)$coefficients

# Creating the coefficient plot
library(ggplot2)
coef_df <- data.frame(
  Predictor = rownames(coefficients),
  Estimate = coefficients[, "Estimate"],
  StdError = coefficients[, "Std. Error"]
)

ggplot(coef_df, aes(x = Estimate, y = reorder(Predictor, Estimate))) +
  geom_point() +
  geom_errorbarh(aes(xmin = Estimate - 1.96 * StdError, xmax = Estimate + 1.96 * StdError), height = 0.1) +
  labs(title = "Coefficient Plot", x = "Estimate", y = "Predictor") +
  theme_minimal()
```



### Explanation of the Coefficient Plot

The coefficient plot is a visual representation of the estimated coefficients for each predictor variable in the regression model. It helps in understanding the strength and direction of the relationships between predictors and the dependent variable.

#### 1. Y-Axis (Predictor):

- The predictor variables are listed along the Y-axis. These are the variables included in the regression model that potentially influence the dependent variable, **AveragScore**.

#### 2. X-Axis (Estimate):

- The X-axis represents the estimated coefficients for each predictor. These coefficients indicate the expected change in the dependent variable for a one-unit change in the predictor, holding all other predictors constant.

#### 3. Point Estimates:

- Each point represents the coefficient estimate for a predictor variable. These estimates indicate how much the dependent variable is expected to change when the predictor changes by one unit.
- For example, if the coefficient for **Education** is 0.083, it means that for each one-unit increase in **Education**, **AveragScore** is expected to increase by 0.083 units, holding all other predictors constant.

#### 4. Error Bars:

- The horizontal lines extending from each point represent the confidence intervals for the coefficient estimates. These intervals provide a range of values within which the true coefficient is expected to fall, with a certain level of confidence (typically 95%).
- Narrower intervals indicate more precise estimates, while wider intervals suggest more uncertainty in the estimate.

## Key Observations:

### 1. Significance and Direction:

- All predictors have positive coefficients, indicating that increases in these predictors are associated with increases in **AveragScore**.
- The magnitude of the coefficients varies, with some predictors having a stronger relationship with the dependent variable than others.

### 2. Relative Importance:

- Predictors such as **EconomicQuality**, **Health**, and **LivingConditions** have higher coefficients, suggesting they have a stronger impact on **AveragScore** compared to others like **SocialCapital** and **NaturalEnvironment**.

### 3. Intercept:

- The intercept term represents the expected value of **AveragScore** when all predictors are zero. In this plot, the intercept is very close to zero, indicating that when all predictors are at their baseline, the **AveragScore** is expected to be nearly zero.

### 4. Precision of Estimates:

- The confidence intervals for most predictors are very narrow, indicating high precision in the estimates. This suggests that the model is well-specified and the data provides strong evidence for the relationships between the predictors and the dependent variable.

## Conclusion:

The coefficient plot provides a clear visualization of the impact of each predictor on the dependent variable. The positive coefficients indicate that higher values of the predictors are associated with higher **AveragScore**. The narrow confidence intervals suggest that these estimates are precise, giving confidence in the reliability of the model.

This plot, combined with the residual plot and other diagnostics, confirms the robustness of the regression model and helps in interpreting the contributions of individual predictors to the overall model.

```
# List of predictor variables
predictors <- names(data_selected)[names(data_selected) != "AveragScore"]

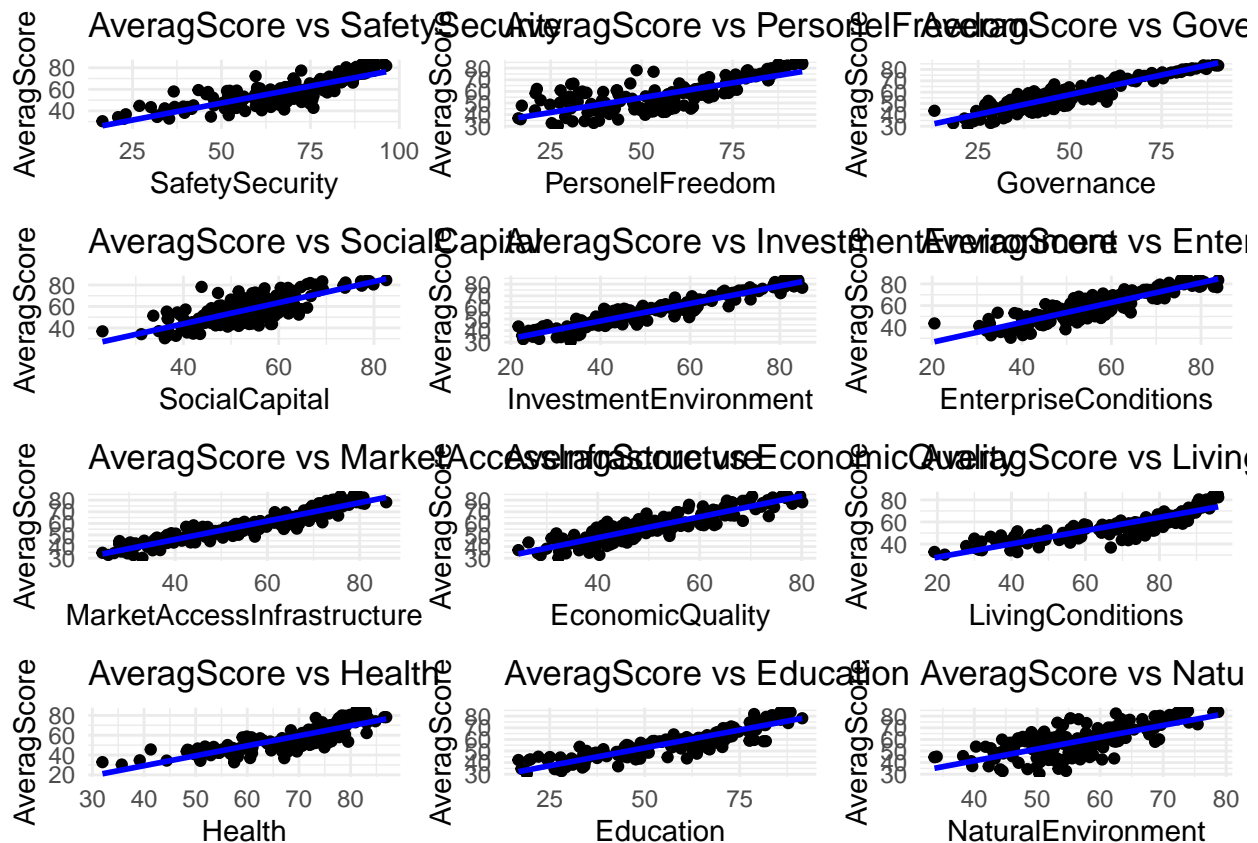
# Function to create scatter plots with regression line
plot_scatter <- function(predictor) {
  ggplot(data_selected, aes_string(x = predictor, y = "AveragScore")) +
    geom_point() +
    geom_smooth(method = "lm", se = FALSE, color = "blue") +
    labs(title = paste("AveragScore vs", predictor), x = predictor, y = "AveragScore") +
    theme_minimal()
}
```

```
# Generate plots for each predictor variable
plots <- lapply(predictors, plot_scatter)
```

```
## Warning: `aes_string()` was deprecated in ggplot2 3.0.0.
## i Please use tidy evaluation idioms with `aes()`.
## i See also `vignette("ggplot2-in-packages")` for more information.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.
```

```
# Arrange plots in a grid
do.call(grid.arrange, c(plots, ncol = 3))
```

```
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
## `geom_smooth()` using formula = 'y ~ x'
```



## Explanation of the Scatter Plots

The scatter plots presented show the relationship between **AveragScore** (y-axis) and various predictor variables (x-axis) in the dataset. Each plot includes a blue fitted line, representing the linear relationship between the two variables.

### 1. **SafetySecurity:**

- A positive linear relationship is observed between **SafetySecurity** and **AveragScore**. As **SafetySecurity** increases, **AveragScore** also tends to increase.

### 2. **PersonelFreedom:**

- There is a strong positive linear relationship between **PersonelFreedom** and **AveragScore**. Higher levels of **PersonelFreedom** correspond to higher **AveragScore**.

### 3. **Governance:**

- The plot shows a positive relationship, indicating that better governance is associated with higher **AveragScore**.

### 4. **SocialCapital:**

- A positive linear trend is evident, suggesting that higher **SocialCapital** leads to higher **AveragScore**.

### 5. **InvestmentEnvironment:**

- A positive relationship is seen, indicating that a better investment environment is associated with higher **AveragScore**.

### 6. **EnterpriseConditions:**

- The plot indicates a positive correlation, with better enterprise conditions being associated with higher **AveragScore**.

### 7. **MarketAccessInfrastructure:**

- There is a positive linear relationship between **MarketAccessInfrastructure** and **AveragScore**.

### 8. **EconomicQuality:**

- A strong positive linear relationship is observed. Higher economic quality corresponds to higher **AveragScore**.

### 9. **LivingConditions:**

- The plot shows a positive linear relationship, indicating that better living conditions are associated with higher **AveragScore**.

### 10. **Health:**

- A positive linear trend is evident. Better health conditions are associated with higher **AveragScore**.

### 11. **Education:**

- There is a strong positive relationship between **Education** and **AveragScore**. Higher levels of education correspond to higher **AveragScore**.

### 12. **NaturalEnvironment:**

- A positive relationship is observed, indicating that a better natural environment is associated with higher **AveragScore**.

## Key Observations:

- **Positive Relationships:**

- All scatter plots exhibit a positive linear relationship between the predictor variables and **AveragScore**. This indicates that improvements in any of these predictors are associated with an increase in the average score.

- **Strength of Relationships:**

- The strength of the relationships varies among the predictors. Variables like **Education**, **PersonelFreedom**, and **EconomicQuality** show particularly strong positive correlations with **AveragScore**.

- **Fitted Lines:**

- The blue fitted lines in each plot provide a clear visualization of the linear trend. The closeness of data points to the fitted lines suggests the degree of fit and the strength of the relationship.

## 5.Conclusion:

These scatter plots visually confirm the positive influence of various predictors on **AveragScore**. The linear relationships suggest that improvements in these areas could lead to higher average scores, supporting the findings from the regression analysis. The fitted lines and the spread of data points around them provide insights into the consistency and strength of these relationships.

## Conclusion and Recommendation:

**Conclusion:** The comprehensive analysis aimed at identifying the key factors influencing the average score in global well-being indices has yielded significant insights. Through various regression models, including simple linear, multiple linear, and stepwise regression methods, we have been able to pinpoint specific predictors that play a crucial role in determining the overall well-being score of countries.

## Key Predictors:

- **Education:** Strongly correlates with higher well-being scores, indicating that better education systems significantly enhance the quality of life.
- **Personal Freedom:** Demonstrates a strong positive impact on well-being, emphasizing the importance of freedom in contributing to overall happiness and satisfaction.
- **Safety and Security:** A critical factor that affects well-being, where safer environments correlate with higher well-being scores.
- **Social Capital:** Shows a significant positive relationship, suggesting that strong social networks and community engagement improve well-being.
- **Economic Quality:** Economic stability and quality directly contribute to higher well-being scores.
- **Health:** Good health and access to healthcare services are vital for higher well-being.
- **Natural Environment:** The quality of the natural environment, including air quality and access to green spaces, positively impacts well-being.

**Model Performance:** The final model, including Education, PersonalFreedom, SafetySecurity, SocialCapital, EconomicQuality, Health, and NaturalEnvironment, demonstrated the best fit. This model had an Adjusted  $R^2$  of 0.983, making it highly robust in explaining the variations in the average well-being score. The model also presented strong values in other metrics, including AIC, BIC, and Mallow's  $C_p$ , showcasing its overall effectiveness.

**Recommendations:** Based on the findings, the following recommendations are proposed to enhance the well-being of populations across countries:

**1. Enhance Educational Systems:**

- Invest in quality education at all levels to ensure that citizens have access to lifelong learning opportunities.
- Implement policies that promote equitable access to education, especially for marginalized and underprivileged communities.

**2. Promote Personal Freedom:**

- Ensure that citizens have the freedom to express themselves, make personal choices, and have control over their lives.
- Protect civil liberties and human rights through robust legal frameworks and enforcement.

**3. Improve Safety and Security:**

- Strengthen law enforcement and public safety measures to create a secure environment for all citizens.
- Implement community policing and other initiatives to build trust between law enforcement and communities.

**4. Strengthen Social Capital:**

- Foster community engagement and social networks through programs that encourage volunteerism and civic participation.
- Support initiatives that promote social cohesion and integration among diverse groups.

**5. Boost Economic Quality:**

- Develop policies that promote economic stability, job creation, and fair wages.
- Support small and medium-sized enterprises (SMEs) to drive economic growth and innovation.

**6. Enhance Healthcare Services:**

- Invest in healthcare infrastructure to provide accessible and affordable health services for all citizens.
- Promote public health initiatives that encourage healthy lifestyles and preventive care.

**7. Preserve the Natural Environment:**

- Implement environmental protection policies to improve air and water quality and conserve natural resources.
- Promote sustainable practices and renewable energy to mitigate the impact of climate change.

**Final Summary:** By focusing on these key areas, countries can significantly improve their well-being scores, leading to happier, healthier, and more prosperous populations. The identified predictors provide a roadmap for policymakers and stakeholders to prioritize and implement effective interventions that enhance the overall quality of life.

## 6. References

- Dunning, J. H. (2002). Determinants of foreign direct investment: Globalization-induced changes and the role of policies. World Investment Report.
- Putnam, R. D. (2000). Bowling alone: The collapse and revival of American community. Simon and Schuster.
- Calderón, C., & Servén, L. (2004). The effects of infrastructure development on growth and income distribution. Policy Research Working Paper Series 3400, The World Bank.
- Hanushek, E. A., & Woessmann, L. (2010). The high cost of low educational performance: The long-run economic impact of improving PISA outcomes. OECD Publishing.
- Sen, A. (1999). Development as freedom. Knopf.
- Wilkinson, R. G., & Pickett, K. (2009). The spirit level: Why more equal societies almost always do better. Allen Lane.
- Kaufmann, D., Kraay, A., & Mastruzzi, M. (2009). Governance matters VIII: Aggregate and individual governance indicators, 1996-2008. World Bank Policy Research Working Paper No. 4978.
- Stiglitz, J. E., Sen, A., & Fitoussi, J. P. (2009). Report by the Commission on the Measurement of Economic Performance and Social Progress. Paris.
- Marmot, M., & Wilkinson, R. G. (Eds.). (2005). Social determinants of health. Oxford University Press.
- McMichael, A. J., Woodruff, R. E., & Hales, S. (2006). Climate change and human health: Present and future risks. *The Lancet*, 367(9513), 859-869.
- United Nations. (2015). Transforming our world: The 2030 agenda for sustainable development.
- Naudé, W. (2010). Entrepreneurship, developing countries, and development economics: New approaches and insights. *Small Business Economics*, 34(1), 1-12.