# MovieLens Project Report

Matheus Assis de Oliveira

November 4, 2024

# Contents

```
## corrplot 0.95 loaded

## Warning: pacote 'ggpubr' foi compilado no R versão 4.4.2

##
## Anexando pacote: 'ggpubr'

## O seguinte objeto é mascarado por 'package:plyr':

##

##     mutate

##
## Anexando pacote: 'dplyr'

## Os seguintes objetos são mascarados por 'package:plyr':

##

##     arrange, count, desc, failwith, id, mutate, rename, summarise,
##     summarize

## Os seguintes objetos são mascarados por 'package:stats':

##

##     filter, lag

## Os seguintes objetos são mascarados por 'package:base':

##

##     intersect, setdiff, setequal, union

## Carregando pacotes exigidos: lattice

## randomForest 4.7-1.2

## Type rfNews() to see new features/changes/bug fixes.

##
## Anexando pacote: 'randomForest'

## O seguinte objeto é mascarado por 'package:dplyr':

##

##     combine

## O seguinte objeto é mascarado por 'package:ggplot2':
```

```
## 
##      margin
```

Packages will be used at the study: tidyverse caret ggplot2 dplyr scales corrplot RColorBrewer ggpubr

# 1   Introduction

The following project is an adaptation of Matheus Assis de Oliveira´s undergraduate thesis in Economics at the Federal University of Mato Grosso do Sul (UFMS), Matheus (2024). The original project aimed to predict the Municipal Human Development Index (IDHM) in Brazil using Machine Learning techniques, particularly Random Forests. The study explored the relationship between the IDHM and various economic and social indicators, such as per capita income, life expectancy, and education levels, to develop a predictive model that could inform policy decisions and development strategies.

The rapid pace of technological development in the digital era has transformed various aspects of society, including the economy, by providing powerful tools for data analysis and forecasting. One of the most promising technologies in this context is Machine Learning (ML), a branch of Artificial Intelligence that enables the identification of patterns in complex datasets and the generation of highly accurate predictions.

In economics, ML is becoming an essential tool for analyzing global variables such as Gross Domestic Product (GDP), population, life expectancy, and average years of schooling. These variables are critical for understanding economic development and informing public policy. However, the application of ML in economic studies is still relatively nascent, offering significant opportunities for innovation and advancement.

Researchers such as Athey and Imbens (Athey and Imbens 2019a) have been at the forefront of efforts to integrate ML into econometrics, highlighting how these techniques can complement traditional methods by handling large datasets and capturing non-linear relationships between variables. For instance, ML can enhance the modeling of economic time series, identify determinants of human development, or predict macroeconomic trends, providing more robust and actionable insights.

This study aims to explore the use of ML techniques to model and predict economic indices using

global variables. By employing models such as Random Forest (Liaw and Wiener 2002a) and comparative methods, the research seeks to demonstrate how data science tools can be integrated into economics to enrich analysis and support strategic decision-making. Furthermore, the study will examine ML's potential to uncover relationships between economic and social variables, contributing to a broader understanding of human and economic development.

Ultimately, this research not only highlights the potential of ML in economics but also provides a practical and accessible example of its application, serving as an introductory guide for economists looking to adopt this technology in their own analytical contexts.

## 2   Literature Review

The application of machine learning (ML) techniques in economics has emerged as an attractive area for academics and professionals. Methods such as Random Forest are widely used to address the complexity and multicollinearity of high-dimensional datasets, making them valuable for analyzing economic and development indicators.

This literature review focuses on prior studies exploring the use of ML in economic research, particularly econometrics, as well as studies that apply predictive methodologies to indicators not strictly economic but with potential applications in economic contexts. Through this review, the aim is to better understand the current state of the field and identify gaps in the existing literature that this work seeks to address.

Hal Ronald Varian, Chief Economist at Google Inc. and Emeritus Professor of Economics at the University of California, Berkeley, has underscored the importance of Big Data and new econometric techniques in his article describing Big Data as a "new trick for economists" (Varian 2014). Although not directly related to predicting economic indicators, Varian's work illustrates the need for economists to become familiar with ML techniques.

Varian explains that ML aims to find a function that provides accurate predictions of "y" based on predictor variables "x." He emphasizes minimizing errors through a loss function and discusses the application of tools such as SQL, Google File System, and BigQuery in data manipulation, alongside models like Random Forest, applied to datasets such as those analyzed by Xavier Sala-i-Martín (Sala-i-Martin 1997).

Another key contributor to the integration of ML in economics is Susan Athey, Professor of Eco-

nomics of Technology at Stanford University. In her 2018 article on the impact of ML on economics, Athey explores the opportunities and challenges of ML applications, including the risks of overfitting and interpretability issues (Athey 2018). She advocates combining economic methods with ML techniques to advance economic research.

Athey defines ML as a field focused on algorithms for datasets, emphasizing supervised learning (prediction and classification) and unsupervised learning (clustering and dimensionality reduction). These methods include k-means clustering, topic modeling, and community detection, which enable economists to analyze large and complex datasets (Athey 2018).

Guido Imbens, a Nobel laureate in Economic Sciences, alongside Athey, further explores ML applications in economics. Their 2019 work provides a comprehensive overview of relevant ML methods for economists, demonstrating their potential for addressing complex economic questions, including human development (Athey and Imbens 2019a).

Studies such as those by Kaur et al. (Kaur et al. 2019) have applied supervised ML techniques, such as logistic regression and Random Forest, to predict Quality of Life indices across nations. These models demonstrated significant predictive power, particularly Random Forest and Support Vector Machines, highlighting ML's utility in economic policy analysis.

Other notable contributions include Sherman et al. (Sherman et al. 2023), who combined satellite imagery with ML to estimate the Human Development Index (HDI) at a high resolution, and Tobaigy et al. (Tobaigy, Alamoudi, and Bafail 2023), who analyzed key determinants of HDI using statistical and ML techniques.

Arumnisaa and Wijayanto (Arumnisaa and Wijayanto 2023a) compared ML methods—Random Forest, Support Vector Machines (SVM), and AdaBoost—for classifying HDI. Their results underline ML's ability to capture non-linear relationships in socio-economic data, aligning with the foundational insights from Athey and Imbens.

In financial markets, ML techniques like Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU) have been used to predict stock prices (Nikou, Mansourfar, and Bagherzadeh 2019; Pesci 2021). These studies reveal ML's superiority over traditional statistical methods in volatile markets.

Internationally, ML has also been applied to human capital analysis in startups (Brigo 2019), energy market price predictions (Vogt 2021), and cryptocurrency price forecasting (Bhattad et al. 2023),

demonstrating its broad applicability across economic domains.

These studies collectively demonstrate ML's potential in capturing complex relationships and generating robust insights. However, they also highlight the need for further research to optimize ML techniques and expand their application to other areas of economic and social development.

In conclusion, the application of ML in economics and human development is a promising and rapidly evolving field. By leveraging advancements in data availability and ML techniques, researchers can enhance the analysis and prediction of critical development indicators such as HDI. Nevertheless, ongoing research is crucial to improve these models' effectiveness and explore their broader applicability.

## 3 Methodology

The Human Development Index (HDI), introduced in 1990 by the United Nations Development Programme (UNDP) under the initiative of Mahbub ul Haq and with contributions from the Indian economist Amartya Sen, represented a significant milestone in quantifying and understanding human development. Published in the *Human Development Report 1990*, the HDI focuses on three fundamental dimensions: health, education, and income (United Nations Development Programme 1990).

While the HDI does not cover all aspects of development, its simplicity and comprehensibility in presenting key indicators have made it an essential tool for assessing human progress globally (Conceição 2022). The HDI evaluates human capabilities by concentrating on three aspects: ensuring a long and healthy life, providing access to knowledge, and ensuring a decent standard of living (United Nations Development Programme, Fundação João Pinheiro, and Instituto de Pesquisa Econômica Aplicada 2013).

In the Brazilian context, the Municipal Human Development Index (IDHM) was developed as a localized adaptation of the global HDI. This adaptation serves to monitor progress at municipal and state levels, integrating specific local data to better reflect the unique characteristics of Brazilian society while maintaining the HDI's core dimensions (United Nations Development Programme, Fundação João Pinheiro, and Instituto de Pesquisa Econômica Aplicada 2013).

This research utilizes Machine Learning (ML) techniques, particularly decision tree-based models such as Random Forests (Liaw and Wiener 2002a), to predict the IDHM. The dataset originates

from the *Atlas of Human Development in Brazil* and is enhanced with global macroeconomic indicators, including population size, Gross Domestic Product (GDP), the proportion of individuals below the poverty line, and infant mortality rates. However, variables like GINI and Theil indices, as well as total population figures, were excluded from the final model due to weak correlations with the unadjusted IDHM ("Atlas Brasil," n.d.).

The HDI is a widely recognized measure for translating human progress into simplified indicators that capture key aspects of human opportunities and capabilities (Varian 2014). Its three core components—health, education, and income—are used to provide a holistic perspective on development:

## 3.1 Long and Healthy Life

The "Long and Healthy Life" dimension focuses on life expectancy as a key indicator. This dimension emphasizes not only increasing life span but also enhancing the quality of life. It reflects the principle that a fulfilling life includes the means to avoid premature death and equitable access to quality healthcare (Kaur et al. 2019).

Life expectancy values are normalized using predefined minimum and maximum benchmarks (25 and 85 years, respectively), following the formula:

$$I = \frac{\text{Observed Value - Minimum Value}}{\text{Maximum Value - Minimum Value}}$$

This normalization facilitates the comparability of longevity across different municipalities and provides a robust basis for policy analysis (Sherman et al. 2023).

## 3.2 Access to Knowledge

The "Access to Knowledge" dimension highlights the transformative role of education in human development. By empowering individuals, fostering autonomy, and building self-esteem, education is framed as both a fundamental right and a strategic tool for informed decision-making about one's future (Athey 2018).

This dimension includes metrics such as adult literacy rates and school enrollment ratios, which are combined into indices ranging from 0 to 1. These indices enable the evaluation of educational

performance across regions, contributing to a nuanced understanding of disparities in access to knowledge (Arumnisaa and Wijayanto 2023a).

## 3.3 Standard of Living

The "Standard of Living" dimension uses income as a proxy for the capacity to fulfill basic needs and achieve a dignified life. Measured through per capita income adjusted for inflation, this indicator reflects the ability of individuals to access essential goods and services such as food, water, and housing (Mullainathan and Spiess 2017).

Normalized income indices consider the diminishing returns of increased income on development, ensuring that the metrics accurately reflect disparities in living standards across municipalities (Jean et al. 2016).

## 3.4 Machine Learning Integration

By integrating ML techniques, this study seeks to enhance the analysis and prediction of the IDHM. Random Forests, known for their robustness in handling complex datasets, are applied to global economic and social variables to model the dimensions of human development (Tobaigy, Alamoudi, and Bafail 2023). This approach demonstrates the potential of ML to provide actionable insights into improving human development globally, particularly at localized levels.

The methodology adopted in this research leverages global economic and social indicators to refine the predictive modeling of the IDHM. By employing ML techniques, this study offers a comprehensive framework for understanding human development, contributing to both academic discourse and practical policy applications (Brigo 2019).

# 4 Data Collection and Preprocessing

The collection and preprocessing of data are fundamental steps in any scientific investigation. In this study, which focuses on the analysis of the Municipal Human Development Index (IDHM), these steps were conducted methodically and with precision.

The primary data source for this study was the Atlas Brasil portal, a reliable and comprehensive repository of socioeconomic data covering all municipalities in Brazil, as well as other federal units and Human Development Units (UDHs). The portal provides data in Excel format (XLSX), which

facilitates importation into data analysis tools such as RStudio ("Atlas Brasil," n.d.).

For this study, IDHM data was collected for all states in Brazil, including the Federal District, organized by territoriality, spanning from 2012 to 2021. To adapt the data for the RStudio platform, the years were renamed with an "X" prefix, resulting in variables such as X2012 through X2021.

Additional indicators collected include literacy rates for individuals aged 15 and older, 18 and older, and 25 and older. Data on the education component of the IDHM, along with its sub-indices for schooling and school attendance, were also retrieved (United Nations Development Programme, Fundação João Pinheiro, and Instituto de Pesquisa Econômica Aplicada 2013).

For the longevity component of the IDHM, life expectancy and infant mortality data were collected. Similarly, the income component included data on per capita income and the percentage of the population living in poverty. To further enhance the breadth of indicators, additional metrics such as total population, GINI index, and Theil index were also included in the dataset (Liaw and Wiener 2002a; Varian 2014).

The preprocessing phase ensured that all data was normalized and structured for analysis, with missing or incomplete entries handled through imputation methods. This comprehensive dataset enables the application of Machine Learning techniques to identify patterns and predict the IDHM effectively.

# 5   Development

After collecting the data, a meticulous data engineering process was undertaken to ensure cleanliness and reliability, which is critical for achieving valid and consistent results. This phase involved identifying and correcting potential inconsistencies in the dataset, such as missing values, errors, and outliers. Different strategies for handling missing values were employed, including mean, median, or mode imputation, depending on the context and nature of the data. Outliers were detected using statistical methods such as the Interquartile Range (IQR) and Z-score techniques and addressed according to the requirements of the study (Irizarry 2019).

Depending on the analysis technique to be applied, the data underwent transformations such as normalization or standardization. These processes ensured that all attributes were on the same scale, making comparisons between variables more consistent and meaningful.

### 5.0.1 Data Engineering Process

To facilitate data engineering, R packages such as `tidyr`, `plyr`, and `readxl` were installed along with their dependencies and loaded using the `library()` function. The first step involved transforming the data into a format suitable for model construction. This included importing data into RStudio and reshaping it using the `gather` function to convert years into a single column, resulting in a table with the variables "Territoriality," "State," "Year" (ranging from 2012 to 2021), and the target variable. The `gsub` function was then used to clean prefixes from year values, and the column was converted to numeric format.

```r
renda_per_capita <- read_excel("renda.per.capita.xlsx")
# Use gather to convert years into a column
renda_per_capita <- gather(renda_per_capita, ano, renda_per_capita,
                           X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
renda_per_capita$ano <- gsub('X', '', renda_per_capita$ano)
# Convert the year column from character to numeric
renda_per_capita <- transform(renda_per_capita, ano = as.numeric(ano))
# Convert the income column to numeric
renda_per_capita <- transform(renda_per_capita, renda_per_capita =
                                  as.numeric(renda_per_capita))


sub_esco_pop <- read_excel("sub.esco.pop.xlsx")
# Use gather to convert years into a column
sub_esco_pop <- gather(sub_esco_pop, ano, sub_esco_pop, X2012:X2021,
                       convert = TRUE)
# Clean the 'X' prefix from the year column
sub_esco_pop$ano <- gsub('X', '', sub_esco_pop$ano)
# Convert the year column from character to numeric
sub_esco_pop <- transform(sub_esco_pop, ano = as.numeric(ano))
# Convert the sub_esco_pop column to numeric
sub_esco_pop <- transform(sub_esco_pop, sub_esco_pop = as.numeric(sub_esco_pop))
```

```r
sub_freq_esco <- read_excel("sub.freq.esco.xlsx")
# Use gather to convert years into a column
sub_freq_esco <- gather(sub_freq_esco, ano, sub_freq_esco, X2012:X2021,
                        convert = TRUE)
# Clean the 'X' prefix from the year column
sub_freq_esco$ano <- gsub('X', '', sub_freq_esco$ano)
# Convert the year column from character to numeric
sub_freq_esco <- transform(sub_freq_esco, ano = as.numeric(ano))
# Convert the sub_freq_esco column to numeric
sub_freq_esco <- transform(sub_freq_esco, sub_freq_esco = as.numeric(sub_freq_esco))


esperança_de_vida <- read_excel("esperança.de.vida.xlsx")
# Use gather to convert years into a column
esperança_de_vida <- gather(esperança_de_vida, ano, esperança_de_vida,
                            X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
esperança_de_vida$ano <- gsub('X', '', esperança_de_vida$ano)
# Convert the year column from character to numeric
esperança_de_vida <- transform(esperança_de_vida, ano = as.numeric(ano))
# Convert the esperança_de_vida column to numeric
esperança_de_vida <- transform(esperança_de_vida, esperança_de_vida =
                                  as.numeric(esperança_de_vida))


porcent_pobres <- read_excel("porcent_pobres.xlsx")
# Use gather to convert years into a column
porcent_pobres <- gather(porcent_pobres, ano, porcent_pobres, X2012:X2021,
                         convert = TRUE)
# Clean the 'X' prefix from the year column
porcent_pobres$ano <- gsub('X', '', porcent_pobres$ano)
# Convert the year column from character to numeric
porcent_pobres <- transform(porcent_pobres, ano = as.numeric(ano))
```

```r
# Convert the porcent_pobres column to numeric
porcent_pobres <- transform(porcent_pobres, porcent_pobres =
                                as.numeric(porcent_pobres))


população_total <- read_excel("população_total.xlsx")
# Use gather to convert years into a column
população_total <- gather(população_total, ano, população_total,
                                X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
população_total$ano <- gsub('X', '', população_total$ano)
# Convert the year column from character to numeric
população_total <- transform(população_total, ano = as.numeric(ano))
# Convert the população_total column to numeric
população_total <- transform(população_total, população_total =
                                as.numeric(população_total))


mortalidade_infantil <- read_excel("mortalidade_infantil.xlsx")
# Use gather to convert years into a column
mortalidade_infantil <- gather(mortalidade_infantil, ano, mortalidade_infantil,
                                X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
mortalidade_infantil$ano <- gsub('X', '', mortalidade_infantil$ano)
# Convert the year column from character to numeric
mortalidade_infantil <- transform(mortalidade_infantil, ano = as.numeric(ano))
# Convert the mortalidade_infantil column to numeric
mortalidade_infantil <- transform(mortalidade_infantil, mortalidade_infantil = as.numeric(morta

media_anos_de_estudo <- read_excel("media_anos_de_estudo.xlsx")
# Use gather to convert years into a column
media_anos_de_estudo <- gather(media_anos_de_estudo, ano, media_anos_de_estudo,
                                X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
```

```r
media_anos_de_estudo$ano <- gsub('X', '', media_anos_de_estudo$ano)
# Convert the year column from character to numeric
media_anos_de_estudo <- transform(media_anos_de_estudo, ano = as.numeric(ano))
# Convert the media_anos_de_estudo column to numeric
media_anos_de_estudo <- transform(media_anos_de_estudo, media_anos_de_estudo = as.numeric(media

indice_gini <- read_excel("indice_gini.xlsx")
# Use gather to convert years into a column
indice_gini <- gather(indice_gini, ano, indice_gini, X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
indice_gini$ano <- gsub('X', '', indice_gini$ano)
# Convert the year column from character to numeric
indice_gini <- transform(indice_gini, ano = as.numeric(ano))
# Convert the indice_gini column to numeric
indice_gini <- transform(indice_gini, indice_gini = as.numeric(indice_gini))

ind_theil_L <- read_excel("ind_theil_L.xlsx")
# Use gather to convert years into a column
ind_theil_L <- gather(ind_theil_L, ano, ind_theil_L, X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
ind_theil_L$ano <- gsub('X', '', ind_theil_L$ano)
# Convert the year column from character to numeric
ind_theil_L <- transform(ind_theil_L, ano = as.numeric(ano))
# Convert the ind_theil_L column to numeric
ind_theil_L <- transform(ind_theil_L, ind_theil_L = as.numeric(ind_theil_L))

analfabetismo_25_anos <- read_excel("analfabetismo_25_anos.xlsx")
# Use gather to convert years into a column
analfabetismo_25_anos <- gather(analfabetismo_25_anos, ano, analfabetismo_25_anos,
                                X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
```

```r
analfabetismo_25_anos$ano <- gsub('X', '', analfabetismo_25_anos$ano)
# Convert the year column from character to numeric
analfabetismo_25_anos <- transform(analfabetismo_25_anos, ano = as.numeric(ano))
# Convert the analfabetismo_25_anos column to numeric
analfabetismo_25_anos <- transform(analfabetismo_25_anos, analfabetismo_25_anos = as.numeric(a


analfabetismo_18_anos <- read_excel("analfabetismo_18_anos.xlsx")
# Use gather to convert years into a column
analfabetismo_18_anos <- gather(analfabetismo_18_anos, ano, analfabetismo_18_anos,
                                X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
analfabetismo_18_anos$ano <- gsub('X', '', analfabetismo_18_anos$ano)
# Convert the year column from character to numeric
analfabetismo_18_anos <- transform(analfabetismo_18_anos, ano = as.numeric(ano))
# Convert the analfabetismo_18_anos column to numeric
analfabetismo_18_anos <- transform(analfabetismo_18_anos, analfabetismo_18_anos = as.numeric(a


analfabetismo_15_anos <- read_excel("analfabetismo_15_anos.xlsx")
# Use gather to convert years into a column
analfabetismo_15_anos <- gather(analfabetismo_15_anos, ano, analfabetismo_15_anos,
                                X2012:X2021, convert = TRUE)
# Clean the 'X' prefix from the year column
analfabetismo_15_anos$ano <- gsub('X', '', analfabetismo_15_anos$ano)
# Convert the year column from character to numeric
analfabetismo_15_anos <- transform(analfabetismo_15_anos, ano = as.numeric(ano))
# Convert the analfabetismo_15_anos column to numeric
analfabetismo_15_anos <- transform(analfabetismo_15_anos, analfabetismo_15_anos = as.numeric(a


IDHM <- read_excel("IDHM.xlsx")
# Use gather to convert years into a column
IDHM <- gather(IDHM, ano, IDHM, X2012:X2021, convert = TRUE)
```

```r
# Clean the 'X' prefix from the year column
IDHM$ano <- gsub('X', '', IDHM$ano)
# Convert the year column from character to numeric
IDHM <- transform(IDHM, ano = as.numeric(ano))
# Convert the IDHM column to numeric
IDHM <- transform(IDHM, IDHM = as.numeric(IDHM))
```

To facilitate the analysis, a unified data frame, named IDHM.AED, was created. This was accomplished using the join function to merge individual data frames containing relevant indicators. The merging was indexed by the fields ano (year) and Territorialidades (territories) to ensure that the data remained organized and correctly aligned. The following code demonstrates this process:

```r
# Initialize the unified data frame with the income data
IDHM.AED <- renda_per_capita


# Join the remaining data frames to create a single unified data frame
IDHM.AED <- join(IDHM.AED, sub_esco_pop, by = c("ano" = "ano", "Territorialidades" = "Territori
IDHM.AED <- join(IDHM.AED, sub_freq_esco, by = c("ano" = "ano", "Territorialidades" = "Territor
IDHM.AED <- join(IDHM.AED, esperança_de_vida, by = c("ano" = "ano", "Territorialidades" = "Ter
IDHM.AED <- join(IDHM.AED, porcent_pobres, by = c("ano" = "ano", "Territorialidades" = "Territo
IDHM.AED <- join(IDHM.AED, população_total, by = c("ano" = "ano", "Territorialidades" = "Territ
IDHM.AED <- join(IDHM.AED, mortalidade_infantil, by = c("ano" = "ano", "Territorialidades" = "T
IDHM.AED <- join(IDHM.AED, media_anos_de_estudo, by = c("ano" = "ano", "Territorialidades" = "T
IDHM.AED <- join(IDHM.AED, indice_gini, by = c("ano" = "ano", "Territorialidades" = "Territoria
IDHM.AED <- join(IDHM.AED, ind_theil_L, by = c("ano" = "ano", "Territorialidades" = "Territoria
IDHM.AED <- join(IDHM.AED, analfabetismo_25_anos, by = c("ano" = "ano", "Territorialidades" = "
IDHM.AED <- join(IDHM.AED, analfabetismo_18_anos, by = c("ano" = "ano", "Territorialidades" = "
IDHM.AED <- join(IDHM.AED, analfabetismo_15_anos, by = c("ano" = "ano", "Territorialidades" = "
IDHM.AED <- join(IDHM.AED, IDHM, by = c("ano" = "ano", "Territorialidades" = "Territorialidades
```

This unified data frame ensures that all indicators relevant to the Municipal Human Development Index (IDHM) are consolidated into a single structure. Each variable maintains its alignment across territories and years, enabling efficient data manipulation and analysis.

Throughout the data engineering process, no errors, null values, or missing values (NA) were found in the unified data frame (IDHM.AED). The creation of the unified data frame also showed no inconsistencies, confirming the reliability and accuracy of the dataset for further analysis.

To validate the dataset, the sapply function was utilized to count the unique values in each column of the data frame. This step aligns with the principles of vectorization and functional programming, as discussed in the section 3.5 of Introduction to Data Science by Rafael Irizarry (Irizarry (2019)). For example, the Territorialidades column displayed 28 unique values, representing Brazil's 26 states, the Federal District, and the national level. The ano column included 10 unique values corresponding to the years from 2012 to 2021. The same validation approach was applied across all other variables, ensuring no duplication or misalignment in the data.

Additionally, the str function was employed to display the structure of the data frame. As outlined in section 2.4.2 of Irizarry's book, this function provides an overview of the data frame's organization and column types. Due to the large size of the data frame, the full structure could not be displayed; however, the output confirmed that all variables were correctly aligned and formatted.

```
AED.df = IDHM.AED
# Using sapply to count unique values for each column
sapply(IDHM.AED, function(x) length(unique(x)))
```

```
##      Territorialidades                   ano       renda_per_capita
##                     28                    10                    280
##           sub_esco_pop          sub_freq_esco       esperança_de_vida
##                    179                   156                    238
##         porcent_pobres        população_total   mortalidade_infantil
##                    261                   280                    255
##   media_anos_de_estudo            indice_gini             ind_theil_L
##                    202                   134                    182
## analfabetismo_25_anos  analfabetismo_18_anos  analfabetismo_15_anos
##                    264                   252                    248
##                   IDHM
##                    135
```

```
# Using str to examine the structure of the data frame
str(IDHM.AED)
```

```
## 'data.frame':    280 obs. of  16 variables:
##  $ Territorialidades   : chr  "Brasil" "Acre" "Alagoas" "Amapá" ...
##  $ ano                 : num  2012 2012 2012 2012 2012 ...
##  $ renda_per_capita    : num  759 517 395 528 559 ...
##  $ sub_esco_pop        : num  0.606 0.59 0.487 0.67 0.613 0.51 0.54 0.765 0.613 0.619 ...
##  $ sub_freq_esco       : num  0.731 0.681 0.645 0.653 0.642 0.639 0.742 0.77 0.735 0.741 .
##  $ esperança_de_vida   : num  74.5 72.5 70 72.8 70.8 ...
##  $ porcent_pobres      : num  11.4 23.8 23.4 18.4 22.2 ...
##  $ população_total     : num  1.98e+08 7.77e+05 3.22e+06 7.21e+05 3.54e+06 ...
##  $ mortalidade_infantil : num  15.8 20.2 26.1 24.3 20.9 ...
##  $ media_anos_de_estudo : num  8.56 7.72 6.8 9.09 8.63 ...
##  $ indice_gini         : num  0.54 0.566 0.503 0.528 0.589 0.563 0.545 0.601 0.489 0.474 .
##  $ ind_theil_L         : num  0.526 0.585 0.447 0.483 0.619 0.571 0.54 0.664 0.411 0.383 .
##  $ analfabetismo_25_anos: num  10.22 18.22 24.22 7.93 9.46 ...
##  $ analfabetismo_18_anos: num  8.75 14.72 20.54 6.37 7.89 ...
##  $ analfabetismo_15_anos: num  8.21 13.48 18.97 5.76 7.22 ...
##  $ IDHM                : num  0.746 0.701 0.651 0.707 0.691 0.678 0.701 0.825 0.758 0.744
```

The results of the sapply and str functions provided confidence in the integrity of the data frame. No missing values (NA) were identified, and all variables demonstrated consistent data entries. The validation ensured that the data frame was error-free and aligned with the analytical requirements.

In addition to technical validation, meticulous documentation was maintained throughout the process to ensure transparency and reproducibility. Each decision was recorded in detail, enabling future researchers to replicate the methodology and validate the findings. This rigorous approach enhances the reliability of the dataset and ensures its readiness for advanced analyses such as predictive modeling.

In conclusion, the data engineering process concluded with a robust verification step, establishing a solid foundation for the subsequent stages of analysis. The careful attention to detail and adherence to methodological rigor ensures that the dataset is both accurate and reliable, serving as a

cornerstone for meaningful and insightful research outcomes.

# 6 Exploratory Data Analysis

Following the data engineering process, a detailed Exploratory Data Analysis (EDA) was conducted to serve as a foundation for understanding and interpreting the underlying complexities of the dataset. This methodology aims to extract intrinsic insights, investigate structural patterns, detect outliers, and test underlying hypotheses that could potentially enhance the Random Forest modeling in subsequent stages.

The initial phase of the exploratory analysis involved a rigorous univariate analysis for each variable in the dataset. This included generating descriptive statistics that provide an understanding of the central characteristics of each variable, such as measures of central tendency (mean, median), dispersion (variance, standard deviation), and extremes (minimum, maximum). Visualizing the distribution of each variable was necessary to understand its shape, identify any deviations from normality, and observe the presence of extreme values.

Visualization techniques, such as scatter plots and correlation heatmaps, were employed to aid in understanding multidimensional relationships. As highlighted by Irizarry (2019), the use of the `dplyr` and `ggplot2` packages enables summarizing and synthesizing data and creating plots and boxplots, respectively. These tools facilitate an intuitive exploration of the dataset's structure.

Outliers, or extreme values, can substantially impact model results. Thus, part of the exploratory analysis focused on detecting and appropriately handling these outliers. Robust methods for identifying these values were implemented, and informed decisions were made regarding whether they should be retained, transformed, or removed based on their nature and impact on the dataset.

## 6.1 Correlation Matrix and Visualization Tools

In the exploratory analysis of the Municipal Human Development Index (IDHM), the packages `corrplot`, `RColorBrewer`, `ggplot2`, and `ggpubr` were utilized. The first step was preparing a dataset to construct a correlation matrix to comprehend the interrelationships among the indicators.

The `corrplot` package (Visualization of a Correlation Matrix) is a powerful tool for visualizing correlation matrices in R, as detailed by Wei and Simko (2021). It provides visual exploratory tools for correlation matrices and supports the automatic reordering of variables, aiding in the detection

of hidden patterns among variables.

The `ggplot2` package, described by Wickham (2016) in *ggplot2: Elegant Graphics for Data Analysis*, is used to create declarative graphics based on *The Grammar of Graphics* Wilkinson (2012). By specifying how variables map to aesthetics and choosing graphical primitives, `ggplot2` takes care of details, enabling the creation of elegant and sophisticated data visualizations.

The `RColorBrewer` package, as explained by Neuwirth (2022), provides color palettes designed by Cynthia Brewer (*ColorBrewer: Color Advice for Maps*, Brewer (2023)). These palettes are effective in representing information clearly and aesthetically.

The `ggpubr` package, discussed by Kassambara (2023), complements `ggplot2` by simplifying the creation of publication-ready plots. While `ggplot2` is flexible and excellent for data visualization, the default plots often require additional formatting. `ggpubr` offers user-friendly functions for creating and customizing `ggplot2` plots of professional quality.

## 6.2 Preparing and Visualizing Correlation Matrices

Initially, non-numeric columns, such as `Territorialidades` and `ano`, were removed, focusing solely on quantitative indicators. Once the data was cleaned, a correlation matrix of the remaining indicators was calculated. Finally, to enhance visualization and interpretation, the `corrplot` package was used to display the matrix (Wei and Simko (2021)).
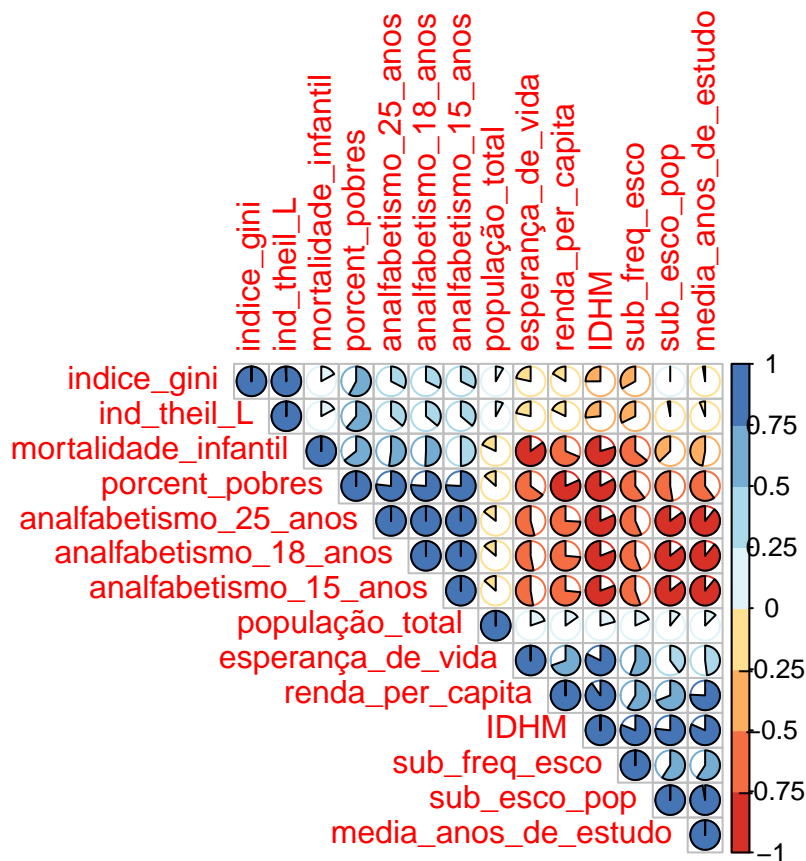
The combination of these tools and methods provided a comprehensive understanding of the dataset's structure and relationships, setting a solid foundation for the predictive modeling stage.

```r
# Create a new data frame excluding rows with missing values
cc = complete.cases(AED.df)
AED.corr = AED.df[cc,]


# Remove non-numeric columns
AED.corr$Territorialidades <- NULL
AED.corr$ano <- NULL


# Compute the correlation matrix
Matrix <- cor(AED.corr)
```

```
# Visualize the correlation matrix with corrplot
corrplot(Matrix,
         type = "upper",        # Display only the upper triangle of the matrix
         order = "hclust",      # Use hierarchical clustering to reorder variables
         method = "pie",        # Use pie charts to represent correlations
         col = brewer.pal(n = 8, name = "RdYlBu")) # Color palette for visualization
```



The creation of a unified data frame, named predict.IDHM, marks a pivotal step in the data preparation process, consolidating all relevant indicators for analysis. This comprehensive dataset integrates key variables, such as per capita income, sub-indicators of education, life expectancy, and demographic metrics, by leveraging the join function to ensure alignment across years and territorial divisions. Following the construction of the data frame, rigorous cleaning processes were undertaken, including the removal of non-numeric columns like Territorialidades and filtering out rows with incomplete data using complete.cases. This ensures the dataset is not only cohesive but also devoid of inconsistencies that might impact subsequent analyses. The final structure of the

data frame was examined using the str function, providing a detailed overview of the variables and confirming its readiness for exploratory and predictive modeling phases. This meticulous preparation enhances the reliability and validity of the findings, laying a solid foundation for further statistical and machine learning applications.

```r
# Create a unified data frame (predict.IDHM) containing all indicators
predic.IDHM <- renda_per_capita


# Join additional indicators using year and territoriality as keys
predic.IDHM <- join(predic.IDHM, sub_esco_pop, by = c("ano" = "ano", "Territorialidades" = "Te
predic.IDHM <- join(predic.IDHM, sub_freq_esco, by = c("ano" = "ano", "Territorialidades" = "T
predic.IDHM <- join(predic.IDHM, esperança_de_vida, by = c("ano" = "ano", "Territorialidades"
predic.IDHM <- join(predic.IDHM, porcent_pobres, by = c("ano" = "ano", "Territorialidades" = "
predic.IDHM <- join(predic.IDHM, população_total, by = c("ano" = "ano", "Territorialidades" =
predic.IDHM <- join(predic.IDHM, mortalidade_infantil, by = c("ano" = "ano", "Territorialides
predic.IDHM <- join(predic.IDHM, media_anos_de_estudo, by = c("ano" = "ano", "Territorialides
predic.IDHM <- join(predic.IDHM, indice_gini, by = c("ano" = "ano", "Territorialidades" = "Ter
predic.IDHM <- join(predic.IDHM, ind_theil_L, by = c("ano" = "ano", "Territorialidades" = "Ter
predic.IDHM <- join(predic.IDHM, analfabetismo_25_anos, by = c("ano" = "ano", "Territorialidade
predic.IDHM <- join(predic.IDHM, analfabetismo_18_anos, by = c("ano" = "ano", "Territorialidade
predic.IDHM <- join(predic.IDHM, analfabetismo_15_anos, by = c("ano" = "ano", "Territorialidade
predic.IDHM <- join(predic.IDHM, IDHM, by = c("ano" = "ano", "Territorialidades" = "Territorial


# Remove non-relevant column (Territorialidades)
predic.IDHM$Territorialidades <- NULL


# Filter rows with complete data
cc = complete.cases(predic.IDHM)
predic.IDHM = predic.IDHM[cc,]


# Display the structure of the final data frame
str(predic.IDHM)
```

```
## 'data.frame':    280 obs. of  15 variables:
##  $ ano                : num  2012 2012 2012 2012 2012 ...
##  $ renda_per_capita   : num  759 517 395 528 559 ...
##  $ sub_esco_pop       : num  0.606 0.59 0.487 0.67 0.613 0.51 0.54 0.765 0.613 0.619 ...
##  $ sub_freq_esco      : num  0.731 0.681 0.645 0.653 0.642 0.639 0.742 0.77 0.735 0.741 .
##  $ esperança_de_vida  : num  74.5 72.5 70 72.8 70.8 ...
##  $ porcent_pobres     : num  11.4 23.8 23.4 18.4 22.2 ...
##  $ população_total    : num  1.98e+08 7.77e+05 3.22e+06 7.21e+05 3.54e+06 ...
##  $ mortalidade_infantil : num  15.8 20.2 26.1 24.3 20.9 ...
##  $ media_anos_de_estudo : num  8.56 7.72 6.8 9.09 8.63 ...
##  $ indice_gini        : num  0.54 0.566 0.503 0.528 0.589 0.563 0.545 0.601 0.489 0.474 .
##  $ ind_theil_L        : num  0.526 0.585 0.447 0.483 0.619 0.571 0.54 0.664 0.411 0.383 .
##  $ analfabetismo_25_anos: num  10.22 18.22 24.22 7.93 9.46 ...
##  $ analfabetismo_18_anos: num  8.75 14.72 20.54 6.37 7.89 ...
##  $ analfabetismo_15_anos: num  8.21 13.48 18.97 5.76 7.22 ...
##  $ IDHM               : num  0.746 0.701 0.651 0.707 0.691 0.678 0.701 0.825 0.758 0.744
```
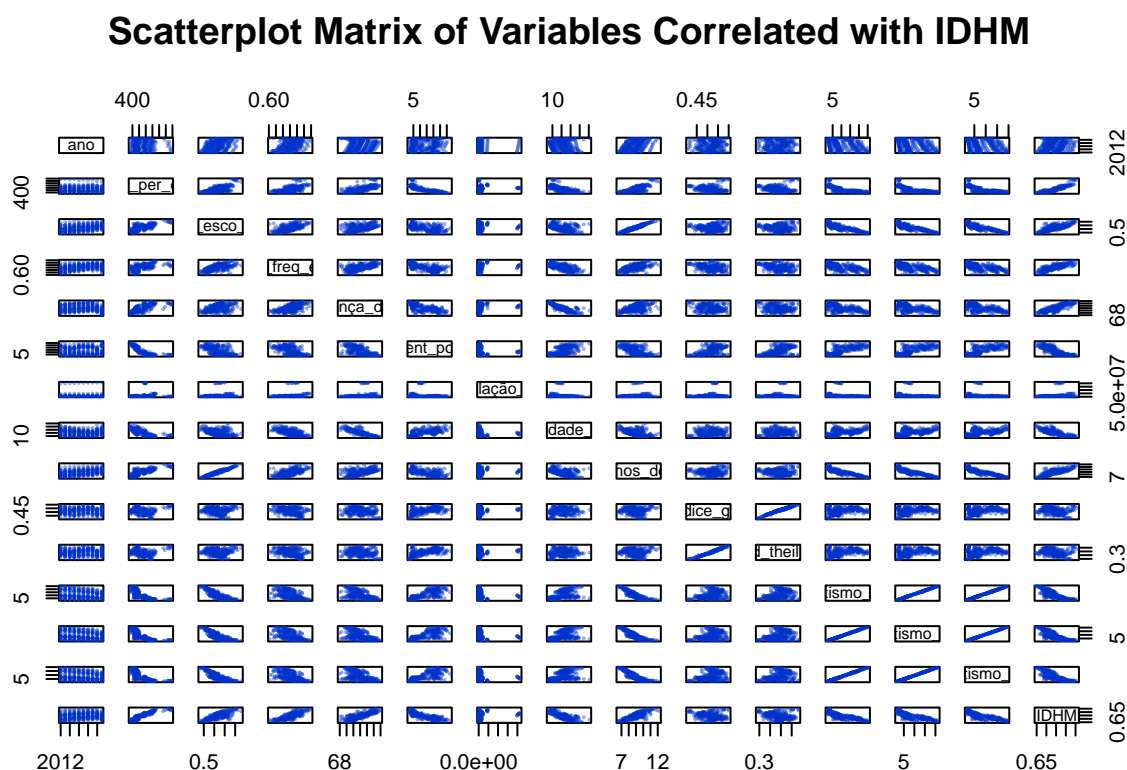
A scatterplot matrix provides an effective visualization tool for understanding the relationships among all variables within the predict.IDHM data frame. This type of visualization is particularly useful for a quick evaluation of interrelations between multiple variables, allowing analysts to identify patterns and dependencies within the dataset. However, a notable drawback of a general scatterplot matrix is the potential for visual overload, especially when the dataset contains numerous variables. This can obscure clear patterns or specific trends, making it challenging to draw actionable insights.

Conversely, plotting the Human Development Index (IDHM) against each individual variable offers a more focused perspective on the specific relationship between the IDHM and other factors. This targeted approach allows for the identification of direct patterns and correlations that are most relevant to understanding the IDHM, avoiding the complexity introduced by examining all possible variable relationships simultaneously.

Both visualization methods were applied in the exploratory data analysis phase to enhance the comprehensiveness of the analysis. The scatterplot matrix served as a broad overview, while individual plots between the IDHM and other variables provided more detailed insights into specific

correlations. The following code demonstrates how the scatterplot matrix was implemented for this analysis:

```
### Plot the correlation as a scatterplot matrix
plot(predic.IDHM, pch=1, cex=0.2, col=rgb(0, 0.2, 0.8, 0.4),
     main = "Scatterplot Matrix of Variables Correlated with IDHM")
```



**Scatterplot Matrix of Variables Correlated with IDHM**

# 7 Model Development and Training

The development and training phase is a fundamental component of the methodological sequence of this investigation. At this stage, the Machine Learning model known as Random Forest was selected due to its notable advantages. Random Forest is particularly suited for handling datasets with complex relationships and is robust against overfitting, a phenomenon where the model learns excessively from the training data and becomes ineffective at predicting new, unseen data.

To prepare for model training, the preprocessed data was split into two subsets: training and testing sets. This common practice is essential for evaluating the performance and generalizability of the

model, providing a robust estimate of its efficacy when applied to new data. The training set is used to fit the model, while the testing set assesses its performance.

Although Random Forest was the primary model of choice due to its advantages, alternative Machine Learning models, such as linear regression, Support Vector Machines (SVM), or neural networks, could also be considered depending on the intrinsic characteristics of the data and the specific problem at hand. Choosing the appropriate model involves careful consideration of trade-offs between model complexity, interpretability, and computational efficiency.

A key element of this phase is hyperparameter tuning, which involves fine-tuning the model's parameters. Adjustments to hyperparameters can significantly impact the model's ability to learn effectively, striking a balance between underfitting (when the model is too simple to capture data complexity) and overfitting (when the model is too complex and overfits the training data, losing generalizability).

Experimenting with different model configurations is an integral part of this process. For instance, altering the structure or depth of the trees in a Random Forest model can yield varying performance outcomes. Such experiments explore a range of scenarios to maximize the model's ability to capture and learn underlying patterns in the data.

The following packages were utilized during the construction and evaluation of the model: plyr, caret, randomForest, and caTools. The caret package offers an integrated approach to training and evaluating Machine Learning models. It simplifies the modeling process by providing a single interface for training various models, coupled with tools for rigorous evaluation, data preprocessing, and visualization Kuhn (2008) . The randomForest package, imported from Breiman's (2001) model, supports classification and regression analyses using a combination of decision trees and random sampling techniques, ensuring more robust and accurate predictions Liaw and Wiener (2002b) . Finally, the caTools package provides versatile utility functions that complement methods like Random Forest and aid data analysis Tuszynski (2021) .

The training process began with an 80/20 split of the dataset into training and testing subsets. This division ensures a balanced approach for training the model and evaluating its performance.

```r
# Split data into 80% training and 20% testing subsets
set.seed(123)
sample.IDHM <- createDataPartition(predic.IDHM$IDHM, p = 0.8, list = FALSE)
```

```
train.IDHM <- predic.IDHM[sample.IDHM, ]
teste.IDHM <- predic.IDHM[-sample.IDHM, ]
```

For the IDHM prediction, the Random Forest model was configured to build multiple decision trees during training. Each tree was generated using a bootstrap sample of the data, and a random subset of features was considered at each node split, determined by the mtry parameter. This randomness serves two main purposes: reducing variance by averaging predictions across trees and decorrelating trees to improve predictive performance.

An initial model (IDHM.model.1) was developed with 500 trees and an arbitrary mtry of 3. The results showed a mean squared residual and an explanation of the variance.

```
# Train a Random Forest model with 500 trees and mtry of 3
IDHM.model.1 <- randomForest(IDHM ~ ., data = train.IDHM, ntree = 500, mtry = 3,
                             importance = TRUE, na.action = na.omit)
print(IDHM.model.1)
```

```
##
## Call:
##  randomForest(formula = IDHM ~ ., data = train.IDHM, ntree = 500,      mtry = 3, importance
##               Type of random forest: regression
##                     Number of trees: 500
## No. of variables tried at each split: 3
##
##          Mean of squared residuals: 5.527219e-05
##                    % Var explained: 97.38
```

To further refine the model, the tuneRF function was used to optimize the mtry parameter. This function iteratively adjusts the value of mtry and evaluates the model using Out-Of-Bag (OOB) error rates. The optimal mtry value found was 3, suggesting the initial configuration was already suitable. Nonetheless, a second model (IDHM.model.2) was developed with an mtry of 4 for educational purposes and comparative analysis.

```
# Optimize mtry using the tuneRF function
mtry_opt <- tuneRF(train.IDHM[-6], train.IDHM$IDHM, ntreeTry = 500,
```

```
                      stepFactor = 1, improve = 0.01, trace = TRUE, plot = FALSE)
```

```
## mtry = 4  OOB error = 1.691716e-05
## Searching left ...
## Searching right ...
```

```
print(mtry_opt)
```

```
##    mtry     OOBError
## 4    4 1.691716e-05
```

```
# Train a second Random Forest model with mtry = 4
IDHM.model.2 <- randomForest(IDHM ~ ., data = train.IDHM, ntree = 500, mtry = 4,
                             importance = TRUE, na.action = na.omit)
print(IDHM.model.2)
```

```
##
## Call:
##  randomForest(formula = IDHM ~ ., data = train.IDHM, ntree = 500,     mtry = 4, importance
##                Type of random forest: regression
##                      Number of trees: 500
## No. of variables tried at each split: 4
##
##          Mean of squared residuals: 5.112506e-05
##                    % Var explained: 97.58
```

The second model achieved a slightly improved variance explanation of. This comparison high-lights the importance of hyperparameter tuning in optimizing model performance. These results underscore the robustness of the Random Forest algorithm and set the stage for the final evaluation phase.

# 8  Model Evaluation and Results

Model evaluation in data science is an indispensable step aimed at assessing the model's effec-tiveness. This phase typically follows model training and has the primary goal of validating the

model's performance against unseen data that was not part of the training process. Beyond merely confirming functionality, model evaluation aids in understanding how predictions are made and helps elucidate the algorithm's practical efficiency.

For this study, the selected Machine Learning model is Random Forest, a robust algorithm capable of managing complex relationships while resisting overfitting. The evaluation of this model was conducted using a separate dataset known as the test set. This partitioning of data into training and testing sets is essential for authentic and unbiased evaluation, providing a reliable estimate of the model's performance on unseen data.

The choice of evaluation metric is critical and directly depends on the type of prediction task. For regression tasks, where the goal is to predict continuous values, the Mean Squared Error (MSE) is commonly used. This metric calculates the average of the squared differences between actual and predicted values, providing a direct numerical measure of model performance.

A high MSE does not necessarily indicate a poor model, nor does a low MSE guarantee a good one. The interpretation of this metric must consider the data scale and the specific context of the problem. Conversely, for classification tasks where the aim is to assign observations to categories, metrics such as accuracy, precision, recall, and F1-score are more relevant. These metrics assess the proportion of correct predictions, the proportion of true positives among all positive predictions, the proportion of true positives among actual positives, and the harmonic mean of precision and recall, respectively.

Another critical aspect of evaluation is balancing performance on the training set and the test set. While a model may perform exceptionally well on training data, poor performance on the test set could indicate overfitting, where the model becomes too tailored to the training data and fails to generalize to new data.

Ultimately, no single metric or evaluation procedure can guarantee a model's effectiveness across all scenarios. The true value of a model lies in its ability to make useful predictions in real-world situations, which may require additional testing and iterative refinement based on feedback and new data. Thus, model evaluation is not merely a step in the development process but an ongoing task that continues even after deployment.

Predictions and Error Metrics In the case of predicting the Municipal Human Development Index (IDHM), the evaluation began by generating predictions using IDHM.model.1 and storing them in

the variable IDHM.predictions.1. For initial insights, predicted IDHM values for six test records were visualized. The same process was repeated for IDHM.model.2, with predictions stored in IDHM.predictions.2.

```
# Generate predictions using Model 1 (mtry = 3)
IDHM.predictions.1 <- predict(IDHM.model.1, teste.IDHM)
head(IDHM.predictions.1)
```

```
##         4         5         6        21        25        26
## 0.7180180 0.7035297 0.6899514 0.7019674 0.7876936 0.7971389
```

```
# Generate predictions using Model 2 (mtry = 4)
IDHM.predictions.2 <- predict(IDHM.model.2, teste.IDHM)
head(IDHM.predictions.2)
```

```
##         4         5         6        21        25        26
## 0.7160924 0.7028467 0.6880818 0.7029156 0.7886690 0.8015258
```

Subsequently, the Root Mean Squared Error (RMSE), a standard regression metric, was calculated to evaluate the difference between observed and predicted values. RMSE provides a square root transformation of MSE, offering an interpretable scale for error measurement.

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \widehat{y}_i)^2}$$

```
# Calculate RMSE for both models
RMSE(IDHM.predictions.1, teste.IDHM$IDHM)
```

```
## [1] 0.008711092
```

```
RMSE(IDHM.predictions.2, teste.IDHM$IDHM)
```

```
## [1] 0.008317511
```

##Variable Importance and Visualization

During the modeling process, variable importance was assessed using the "Percentage Increase in Mean Squared Error" (%IncMSE) and "Increase in Node Purity" (IncNodePurity). These metrics measure how a variable influences model accuracy and node homogeneity, respectively. Key

variables such as per capita income, school attendance subindex, and life expectancy showed high importance across both metrics. Conversely, the year variable had minimal weight, indicating lower relevance in predicting IDHM.

```
# Assess variable importance
importance(IDHM.model.1)
```

```
##                         %IncMSE IncNodePurity
## ano                    4.889520   0.002267137
## renda_per_capita      20.466047   0.106702537
## sub_esco_pop          17.031738   0.019183027
## sub_freq_esco         21.193368   0.029430968
## esperança_de_vida     19.738248   0.040186554
## porcent_pobres        12.827685   0.054979721
## população_total       13.679770   0.005553883
## mortalidade_infantil  16.808312   0.055097987
## media_anos_de_estudo  14.732096   0.027432311
## indice_gini           11.003634   0.003372158
## ind_theil_L            9.703659   0.003630963
## analfabetismo_25_anos 13.714685   0.045379432
## analfabetismo_18_anos 11.726397   0.045213473
## analfabetismo_15_anos 11.533864   0.033028870
```

```
importance(IDHM.model.2)
```

```
##                         %IncMSE IncNodePurity
## ano                    7.632053   0.001784957
## renda_per_capita      20.835682   0.127887976
## sub_esco_pop          17.645840   0.016095339
## sub_freq_esco         25.414882   0.023907318
## esperança_de_vida     20.741278   0.048109606
## porcent_pobres        12.733587   0.063838444
## população_total       11.103371   0.003869649
## mortalidade_infantil  16.647434   0.054706016
```

```
## media_anos_de_estudo  17.754598   0.027180729
## indice_gini            9.131258   0.002798168
## ind_theil_L            8.293696   0.002733618
## analfabetismo_25_anos 12.189939   0.040030954
## analfabetismo_18_anos 11.647392   0.035036142
## analfabetismo_15_anos  9.923677   0.023800294
```
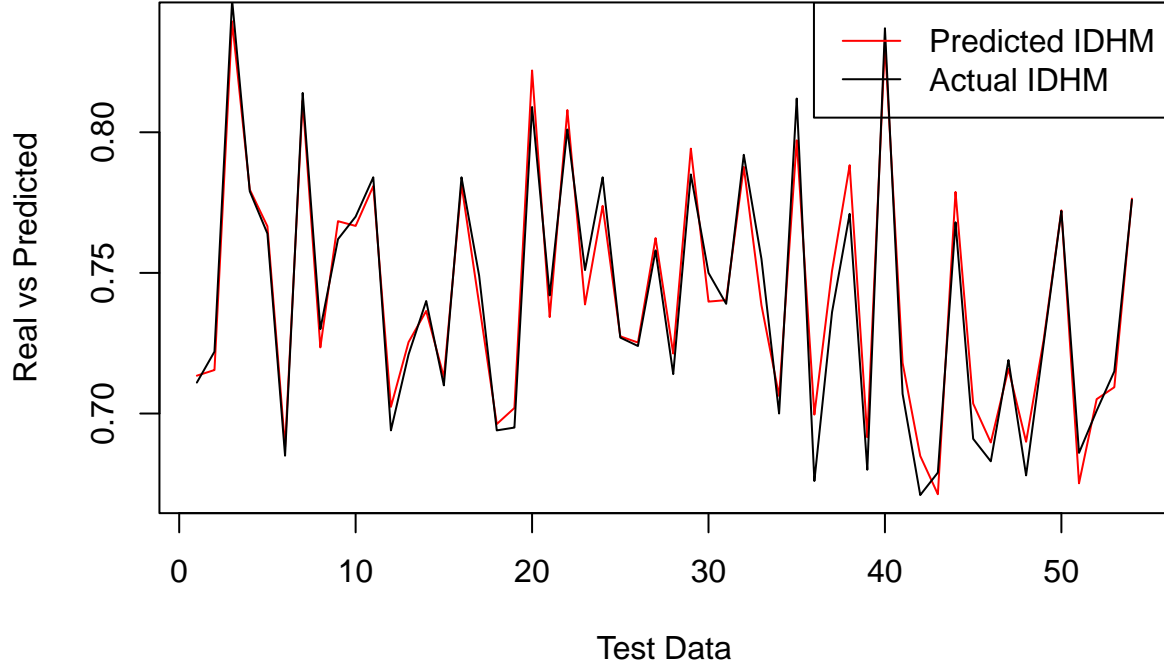
To finalize the evaluation, predicted and actual IDHM values were aligned in a single data frame, enabling a detailed comparison. A new column, diff, was created to calculate the difference between real and predicted values for each observation, providing granular insights into model performance.

```r
# Create a data frame with actual and predicted values, and compute the difference
IDHM.predictions.df <- data.frame(IDHM.predictions.1)
IDHM.predictions.df <- merge(teste.IDHM, IDHM.predictions.df, by.x = 0, by.y = 0)
IDHM.predictions.df$diff <- IDHM.predictions.df$IDHM - IDHM.predictions.df$IDHM.predictions.1
```

Finally, the evaluation concluded with a visualization comparing real and predicted IDHM values. This graphical representation highlights the model's predictive accuracy across the test set.

```r
# Plot actual vs. predicted values
plot(IDHM.predictions.df$IDHM.predictions.1, type = "l", col = "red",
     xlab = "Test Data", ylab = "Real vs Predicted", main = "IDHM Predictions")
lines(IDHM.predictions.df$IDHM, col = "black")
legend("topright", legend = c("Predicted IDHM", "Actual IDHM"), col = c("red", "black"), lty =
```

## IDHM Predictions



In summary, the Random Forest model demonstrated robust performance in predicting IDHM, achieving an RMSE of 0.00817621 in the second configuration. The evaluation underscores the importance of careful tuning and comparison of model configurations to optimize predictive performance(Breiman 2001a), (Liaw and Wiener 2002b), (Kuhn 2008).

## 9 CONCLUSION

This study highlighted the increasing adoption of machine learning (ML) in the field of economics, underscoring its appeal among renowned economists, chief economists at major corporations, and even Nobel laureates in Economic Sciences. Figures such as Professor Guido Imbens and Google's Chief Economist Hal Varian exemplify the growing recognition of ML as a vital tool for economic analysis and decision-making (Athey and Imbens 2019b; Varian 2014).

One of the central challenges in applying ML models lies in interpreting the results and effectively communicating them to diverse audiences. This process requires interdisciplinary skills that extend beyond the technical aspects of ML, encompassing elements of statistics, computer science, data

visualization, and communication (Irizarry 2019).

A key component of interpreting ML results is evaluating variable importance within the model. For Random Forests, this is often achieved using permutation importance, where the model's accuracy is assessed after randomly shuffling a variable's values while keeping others constant. Variables that cause a significant drop in accuracy are deemed more important. However, this method should be used cautiously, as the importance of variables can be influenced by factors such as variable correlations and measurement scales. Complementary techniques like principal component analysis or correlation heatmaps can provide additional insights to enhance interpretability (Liaw and Wiener 2002a).

The Random Forest model demonstrated excellent fit during training, explaining the variance without overfitting, as it did not reach a perfect explanation of 100%. This result indicates that the model achieved a low mean squared error, demonstrating minimal prediction error for the Human Development Index (HDI), which ranges from 0 to 1 (Breiman 2001b).

Further evaluation showed an average prediction error of, highlighting that, while the model performs well, a small margin of variation in predictions remains. This aligns with findings in other studies, such as Arumnisaa and Wijayanto (2023b), where Random Forest models achieved an accuracy of 85.23% in classifying HDI for districts and cities in Indonesia, outperforming Support Vector Machines (84.61%) and AdaBoost (80.36%). These comparisons confirm the reliability and high efficiency of the Random Forest model employed in this study.

This research contributes to the dissemination of novel tools in economics, providing a step-by-step practical demonstration of ML's application. It serves as a resource for newcomers aiming to develop their first ML models within the field of economics. Additionally, this work builds upon prior studies where ML was applied not only to predict indicators like HDI but also to analyze economic growth, forecast energy prices in Germany (Vogt 2021), and predict stock prices ( 2021). These examples underscore ML's versatility in modeling and forecasting economic variables, positioning it as an invaluable tool for economists to learn and replicate.

# References

Arumnisaa, R. I., and A. W. Wijayanto. 2023a. "Comparison of Ensemble Learning Method: Random Forest, Support Vector Machine, AdaBoost for Classification Human Development

Index (HDI)." *Sistemasi: Jurnal Sistem Informasi* 12 (1): 206–18.

———. 2023b. "Comparison of Ensemble Learning Method: Random Forest, Support Vector Machine, AdaBoost for Classification Human Development Index (HDI)." *Sistemasi: Jurnal Sistem Informasi* 12 (1): 206–18.

Athey, Susan. 2018. "The Impact of Machine Learning on Economics." *National Bureau of Economic Research.* https://www.nber.org/books-and-chapters/economics-artificial-intelligence-agenda/impactmachine-learning-economics.

Athey, Susan, and Guido W. Imbens. 2019a. "Machine Learning Methods That Economists Should Know About." *Annual Review of Economics* 11 (1): 685–725.

———. 2019b. "Machine Learning Methods That Economists Should Know About." *Annual Review of Economics* 11 (August): 685–725. https://doi.org/10.1146/annurev-economics-080217-053433.

"Atlas Brasil." n.d. Online. http://www.atlasbrasil.org.br/acervo/atlas.

Bhattad, S. et al. 2023. "Review of Machine Learning Techniques for Cryptocurrency Price Prediction." EasyChair Preprint. https://easychair.org/publications/preprint_open/t5fX.

Breiman, Leo. 2001b. "Random Forests." *Machine Learning* 45 (1): 5–32. https://doi.org/10.1023/A:1010933404324.

———. 2001a. "Random Forests." *Machine Learning* 45 (1): 5–32. https://doi.org/10.1023/A:1010933404324.

Brewer, Cynthia. 2023. "ColorBrewer: Color Advice for Maps."

Brigo, Francesco. 2019. "Applicazione Di Tecniche Di Machine Learning Per l'analisi Del Ruolo Del Capitale Umano Nelle Startup."

Conceição, P. 2022. *Human Development Report.* Lanham: Bernan Press.

Irizarry, Rafael A. 2019. *Introduction to Data Science.* Boca Raton, FL: CRC Press.

Jean, Neal et al. 2016. "Combining Satellite Imagery and Machine Learning to Predict Poverty." *Science* 353 (6301): 790–94.

Kassambara, Alboukadel. 2023. *Ggpubr: 'Ggplot2' Based Publication Ready Plots.*

Kaur, M. et al. 2019. "Supervised Machine-Learning Predictive Analytics for National Quality of Life Scoring." *Applied Sciences* 9 (8): 1613.

Kuhn, Max. 2008. "Building Predictive Models in r Using the Caret Package." *Journal of Statistical Software* 28 (5): 1–26. https://doi.org/10.18637/jss.v028.i05.

Liaw, Andy, and Matthew Wiener. 2002a. "Classification and Regression by randomForest." *R*

*News* 2 (3): 18–22. https://CRAN.R-project.org/doc/Rnews/.

———. 2002b. "Classification and Regression by randomForest." *R News* 2 (3): 18–22. https://CRAN.R-project.org/doc/Rnews/.

Matheus. 2024. "Aplicação Do Machine Learning Na Previsão de Índices Econômicos: O IDHM Com o Modelo Random Forest de 2012 à 2021."

Mullainathan, Sendhil, and Jann Spiess. 2017. "Machine Learning: An Applied Econometric Approach." *The Journal of Economic Perspectives* 31 (2): 87–106.

Neuwirth, Erich. 2022. *RColorBrewer: ColorBrewer Palettes.*

Nikou, Mahla, Gholamreza Mansourfar, and Jamshid Bagherzadeh. 2019. "Stock Price Prediction Using Deep Learning Algorithm and Its Comparison with Machine Learning Algorithms." *Intelligent Systems in Accounting, Finance and Management.*

Pesci, P. 2021. "Previsione Del Prezzo Delle Azioni Di s&p Con Reti Neurali LSTM e GRU."

Sala-i-Martin, Xavier X. 1997. "I Just Ran Two Million Regressions." *The American Economic Review* 87 (2): 178–83.

Sherman, L. et al. 2023. "Global High-Resolution Estimates of the United Nations Human Development Index Using Satellite Imagery and Machine-Learning."

Tobaigy, Faisal, Mohammed Alamoudi, and Omar Bafail. 2023. "Human Development Index: Determining and Ranking the Significant Factors." *International Journal of Engineering Research & Technology* 12 (3).

Tuszynski, Jarek. 2021. "caTools: Tools: Moving Window Statistics, GIF, Base64, ROC AUC, Etc." https://cran.r-project.org/web/packages/caTools/index.html.

United Nations Development Programme. 1990. "Human Development Report 1990: Concept and Measurement of Human Development." New York: United Nations Development Programme.

United Nations Development Programme, Fundação João Pinheiro, and Instituto de Pesquisa Econômica Aplicada. 2013. *O Índice de Desenvolvimento Humano Municipal Brasileiro.* Brasília, Distrito Federal, Brazil: PNUD.

Varian, Hal R. 2014. "Big Data: New Tricks for Econometrics." *Journal of Economic Perspectives* 28 (2): 3–28.

Vogt, Jan. 2021. "Vorhersage von Aktienkursbewegungen Der Energiebranche Mithilfe Maschinellen Lernens Und Stimmungserkennung von Beiträgen Aus Sozialen Medien."

Wei, Taiyun, and Viliam Simko. 2021. *Corrplot: Visualization of a Correlation Matrix.*

Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis.* Springer-Verlag New York.

Wilkinson, Leland. 2012. *The Grammar of Graphics.* Springer.

, . . 2021. " ." 6: 21–26.