# "Predicting the Age of Abalone Using Regression Models

## A Comparative Study of Linear, Ridge, and Lasso Regression"

**MUHAMMAD MUBASHAR SHAHZAD**

**Registration No. SM3600012**

**Scientific and Data Intensive Computing (SDIC)**

University of Trieste

# Introduction

- Objective: To predict the age of abalone using physical measurements.

- Dataset: Abalone dataset with 4177 instances and 8 features.

- Target Variable: Rings (predicting Age by adding 1.5).

# Dataset Overview

- Features:
- Sex (Categorical)
- Length (Continuous)
- Diameter (Continuous)
- Height (Continuous)
- Whole weight (Continuous)
- Shucked weight (Continuous)
- Viscera weight (Continuous)
- Shell weight (Continuous)
- Target: Rings (Integer, Age = Rings + 1.5)

# Data Loading and Initial Exploration

- Step 1: Loading the dataset using Pandas.

- Code Snippet: `pd.read_csv('abalone.csv')`

- Initial DataFrame: Show the first few rows of the dataset.

- Handling Missing Values: Mention that there are no missing values in the dataset.
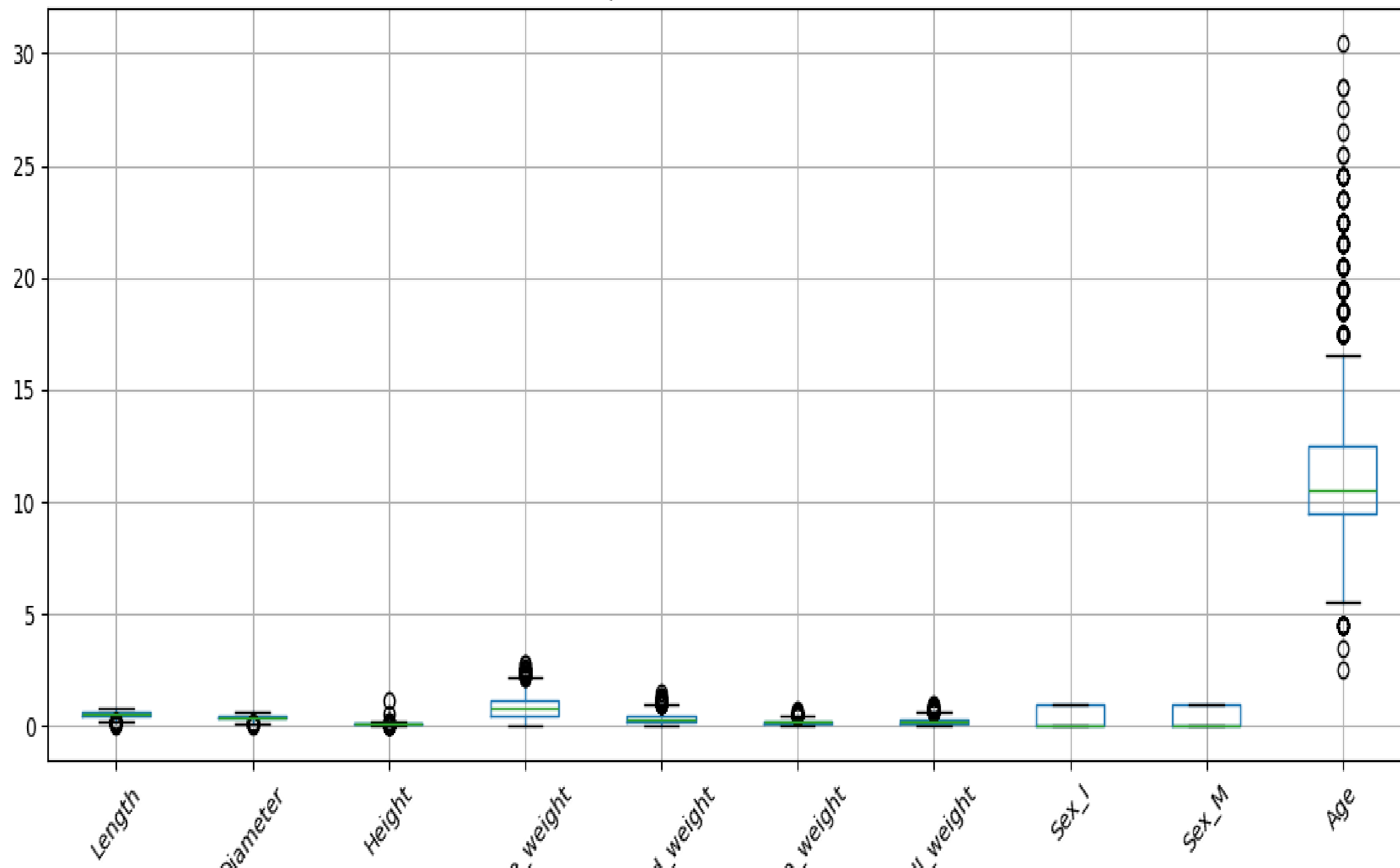
# Data Preprocessing

- Step 2: Converting Categorical Variable

- Sex column converted to numerical using one-hot encoding.

- Step 3: Creating Age Column
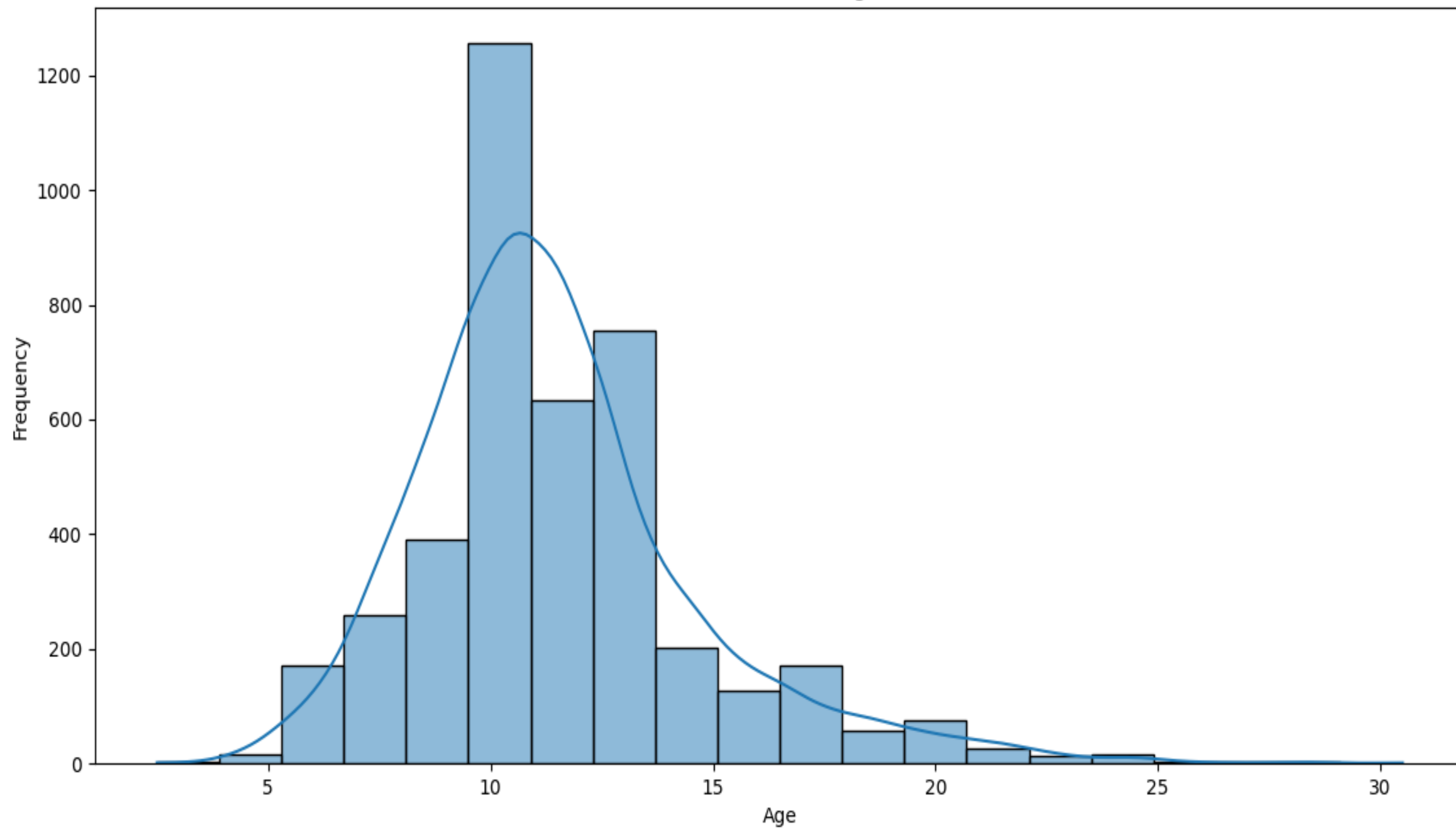
- Adding 1.5 to the Rings column to get the age.

# Exploratory Data Analysis (EDA)

➡ Step 4: Visualizations

➡ Histograms of numerical features.

➡ Box plots to check for outliers.

➡ Step 5: Correlation Analysis

➡ Heatmap of correlation matrix to understand
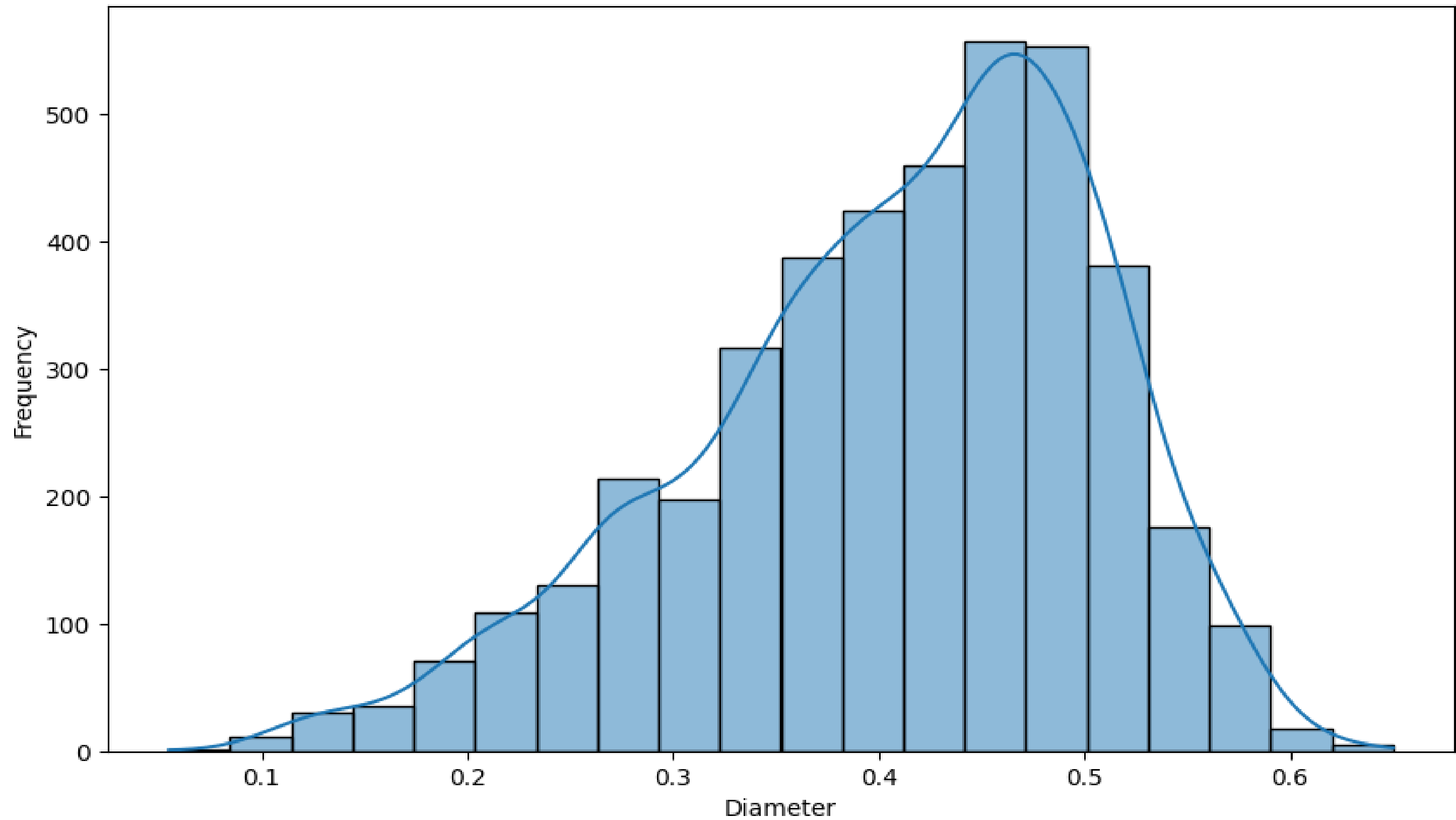
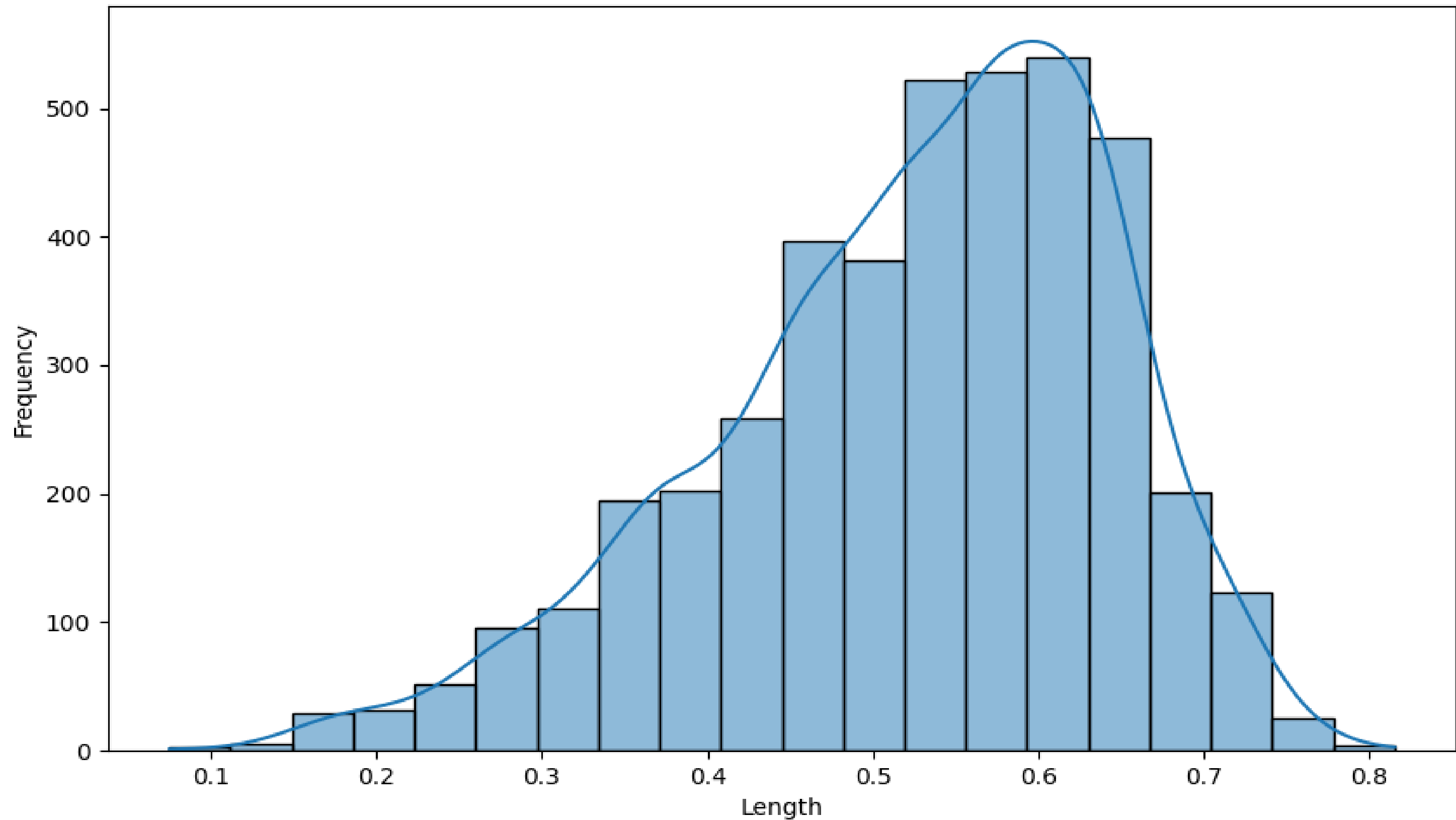relationships between features and target.

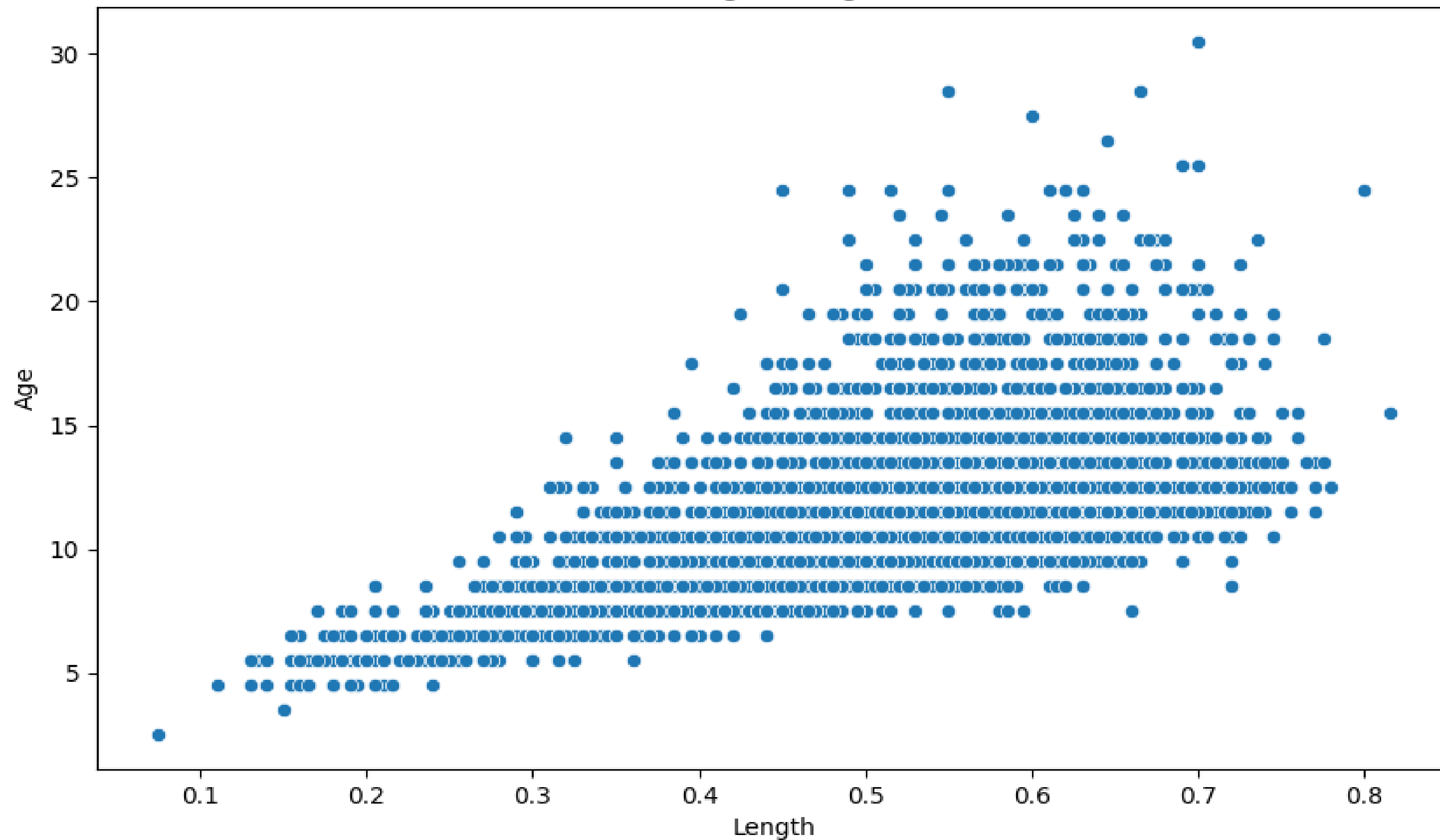Boxplot of Numerical Features

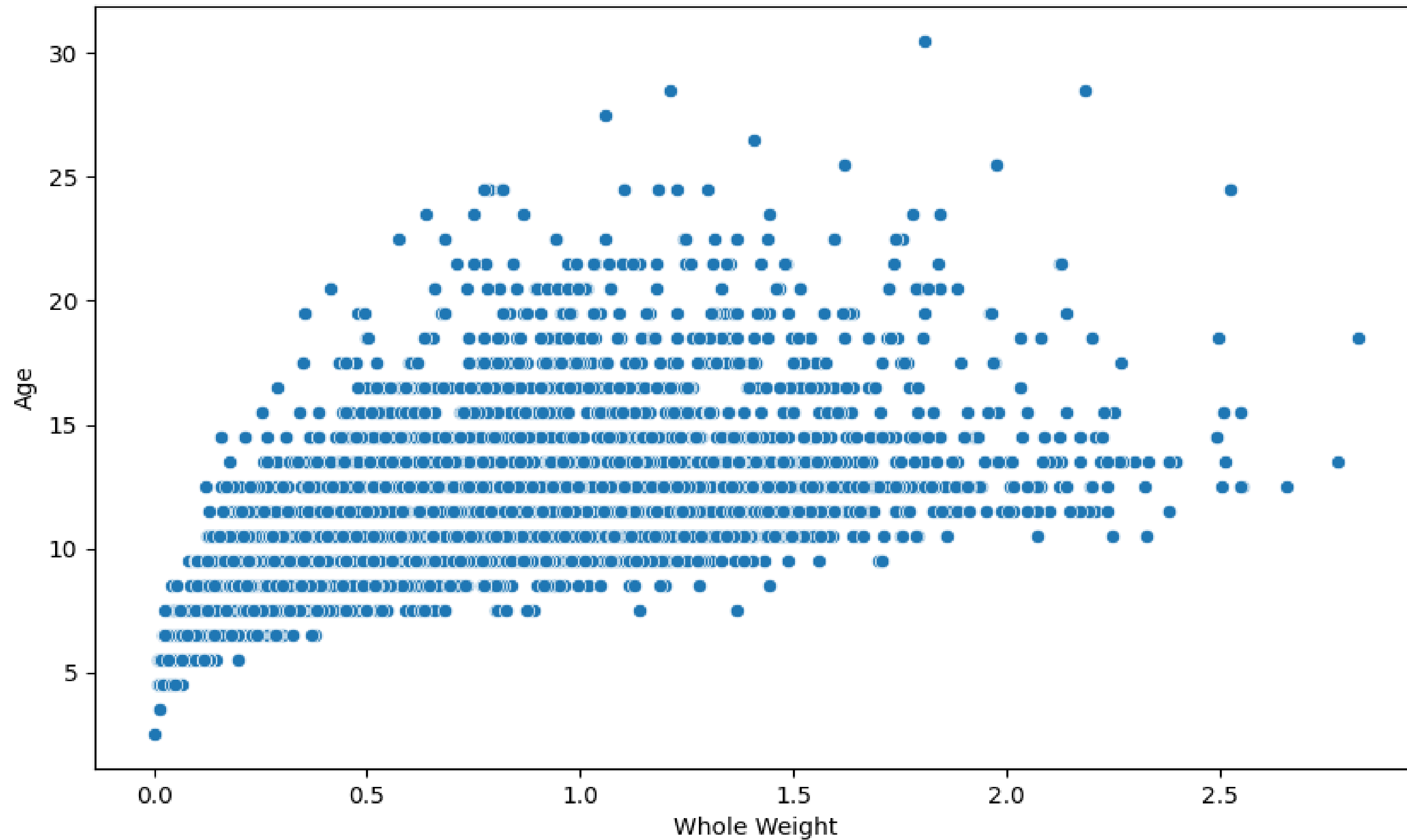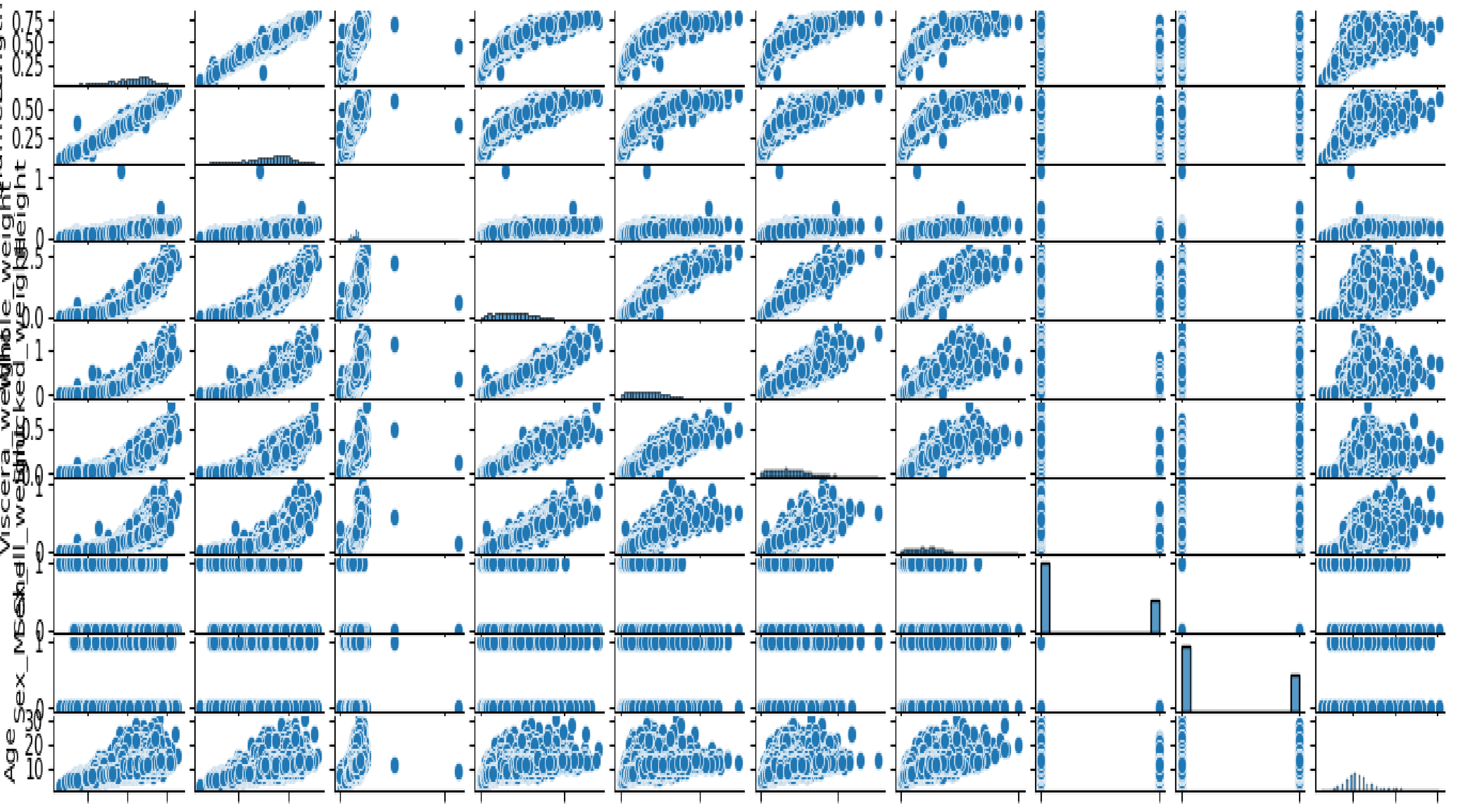Distribution of Age

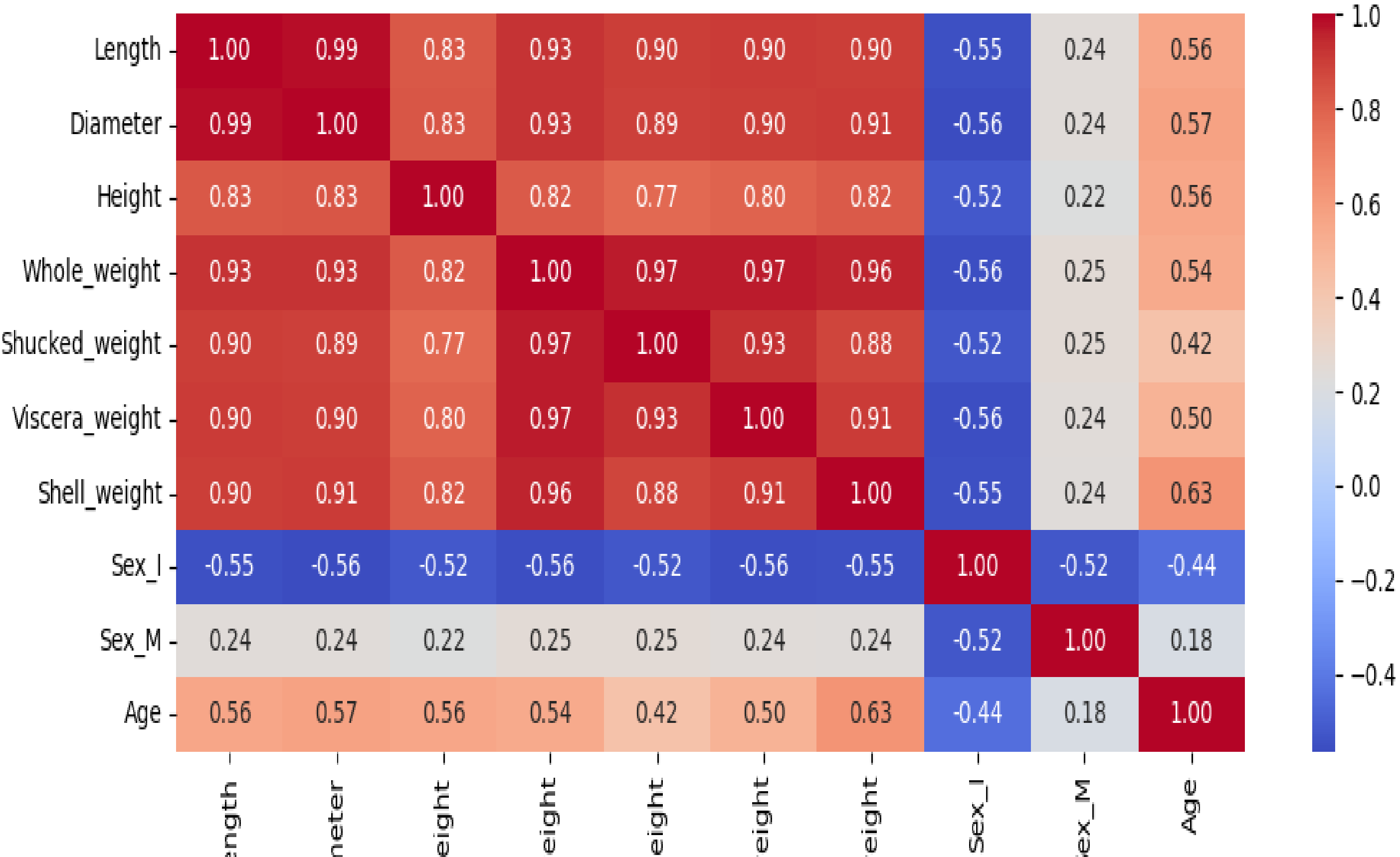Distribution of Diameter

Distribution of Length

Length vs Age

Whole Weight vs Age

Correlation Matrix

# Feature Scaling and Train-Test Split

- Step 6: Scaling the Features

- Normalization using StandardScaler.

- Step 7: Splitting the Data

- Train-Test split (80-20) using `train_test_split`.

# Model Training - Linear Regression

- Model: Linear Regression

- Training: Fit the model on training data.

- Evaluation: MSE and R-squared on test data.

- Results:

- MSE: 4.8912

- -R2: 0.5482

# Model Training - Ridge Regression

- Model: Ridge Regression

- Training: Fit the model on training data.

- Evaluation: MSE and R-squared on test data.

- Results:

- MSE: 4.8911

- R2: 0.5482

# Model Training - Lasso Regression

- Model: Lasso Regression

- Training: Fit the model on training data.

- Evaluation: MSE and R-squared on test data.

- Results:

- MSE: 7.6826

- R2: 0.2903

# Model Comparison

- Linear Regression:
- MSE: 4.8912
- R2: 0.5482
- Lasso Regression:
- MSE: 7.6826
- R2: 0.2903
- Ridge Regression:
- MSE: 4.8911
- R2: 0.5482

# Conclusion

- Summary: Linear and Ridge Regression perform similarly and better than Lasso Regression.

- Recommendation: Choose Ridge Regression if multicollinearity is a concern; otherwise, Linear Regression is also a good choice.

- Future Work: Consider other regression techniques or feature engineering for improved performance.

# Questions & Answers

Open the floor for any questions from the audience.