



DSEM 6150 ASSIGNMENT 6

XGBOOST ALGORITHM

P R E S E N T E D B Y : R U C H I K A P A D I W A L A



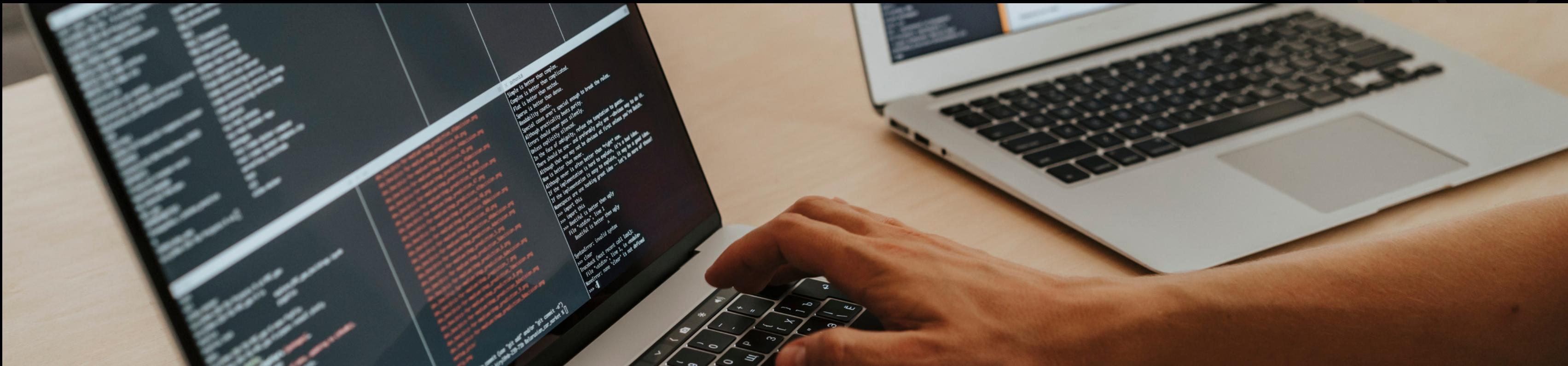
INTRODUCTION

XGBoost(Extreme Gradient Boosting) is a popular machine learning algorithm that is widely used for solving regression and classification problems.

It is an implementation of the gradient boosting algorithm, which builds an ensemble of weak decision trees and combines them to make a strong prediction.

XGBoost is known for its scalability, speed, and high predictive accuracy, making it a popular choice for many data science applications.





PROBLEM STATEMENT

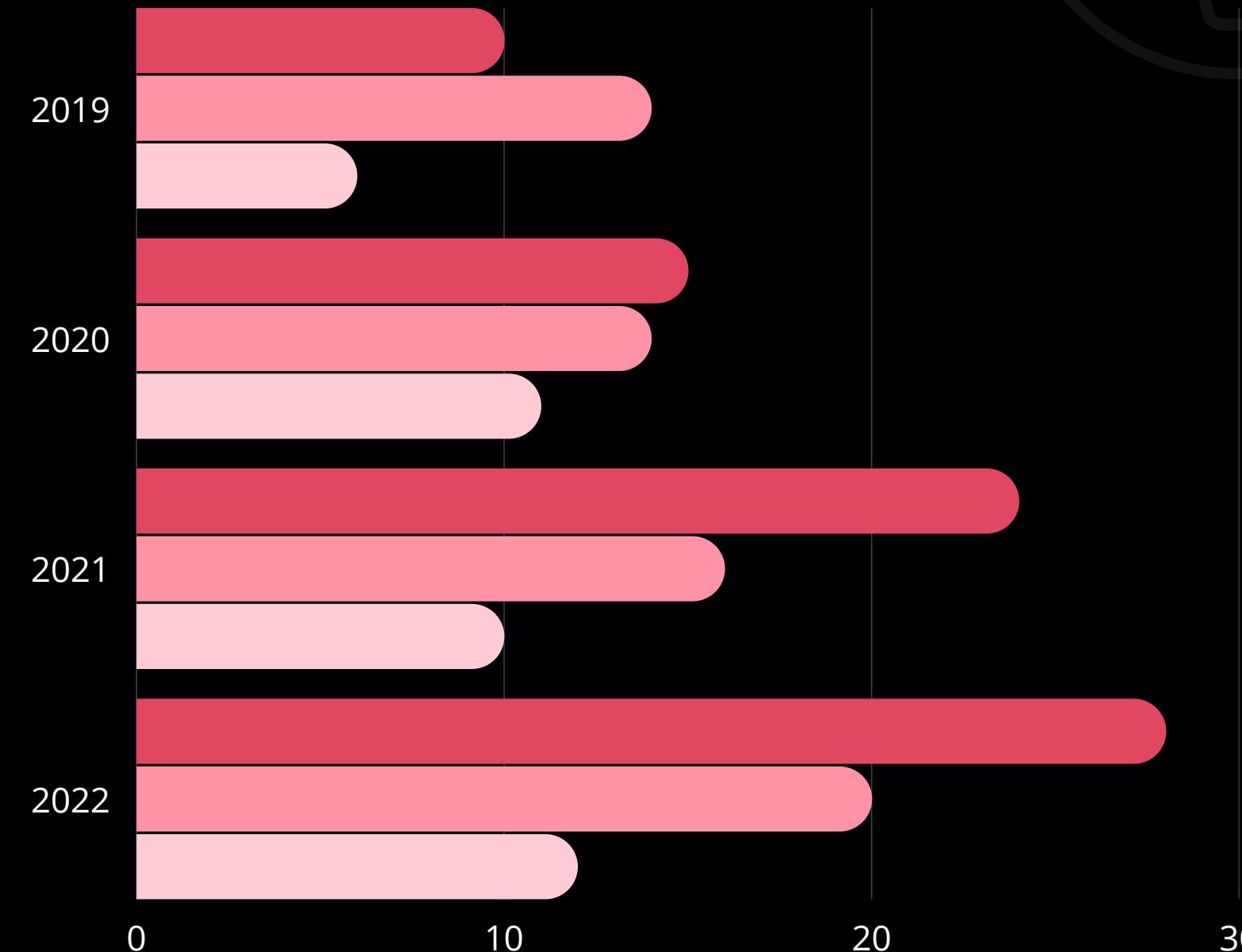
Applying XGBoost algorithm on the Boston housing dataset to predict the median value. Do the usual train/test split and then compare the result with the result of applying linear regression. You can use a suitable metric like RMSE.

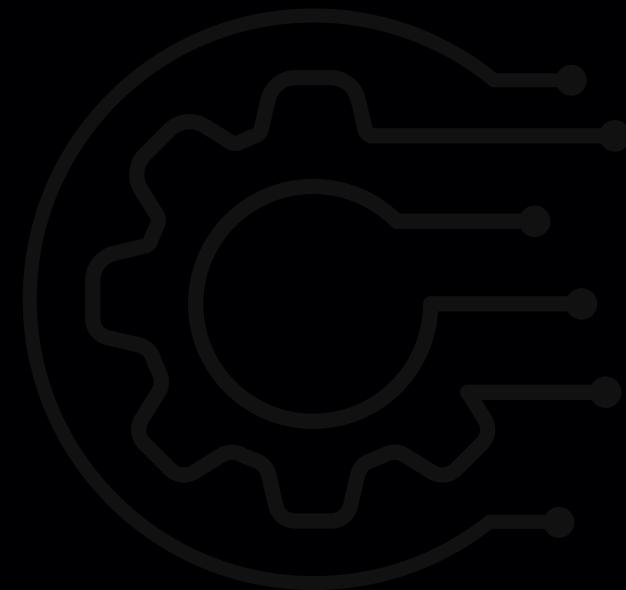


DATA DESCRIPTION

Description

The Boston Housing Dataset is derived from information collected by the U.S. Census Service concerning housing in the area of Boston MA. The following describes the dataset columns:





01

CRIM

per capita crime rate by town

02

ZN

proportion of residential land zoned for lots over
25,000 sq.ft.

03

INDUS

proportion of non-retail business acres per town.

04

CHAS

Charles River dummy variable (1 if tract bounds
river; 0 otherwise)

05

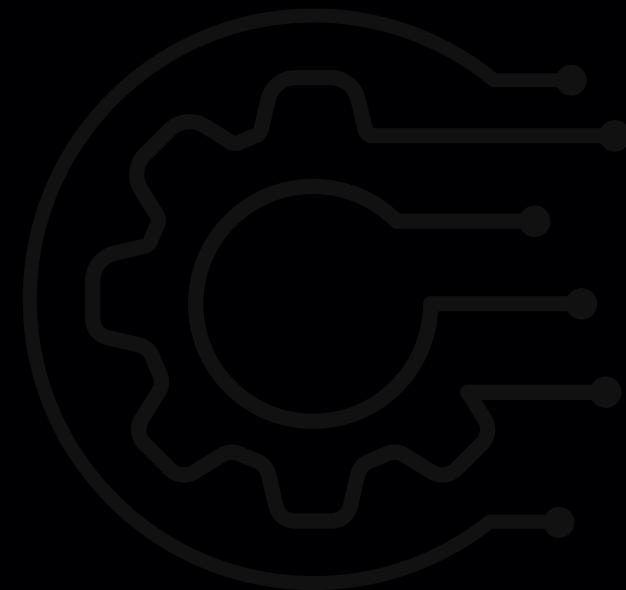
NOX

nitric oxides concentration (parts per 10 million)

06

RM

average number of rooms per dwelling



07

AGE

proportion of owner-occupied units built prior to 1940

08

DIS

weighted distances to five Boston employment centres

09

RAD

index of accessibility to radial highways

10

PTRATIO

full-value property-tax rate per \$10,000

11

B

$1000(B_k - 0.63)^2$ where B_k is the proportion of blacks by town

12

LSTAT

% lower status of the population



13

MV

Median value of owner-occupied homes in \$1000's



TECHNICAL APPROACH



Data preprocessing: Cleaning and transforming data to prepare it for modeling



Splitting Target Variable and Independent Variable

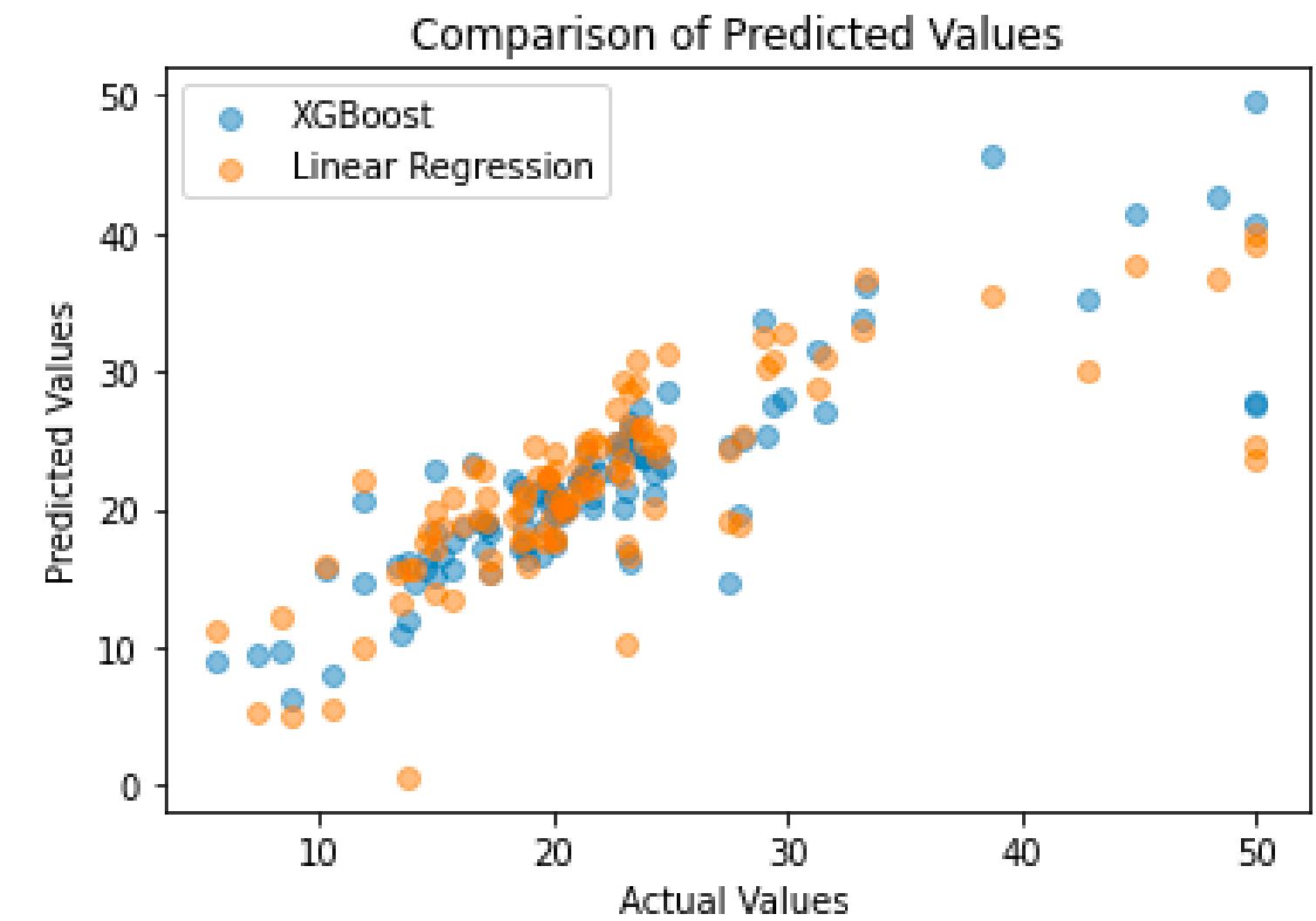
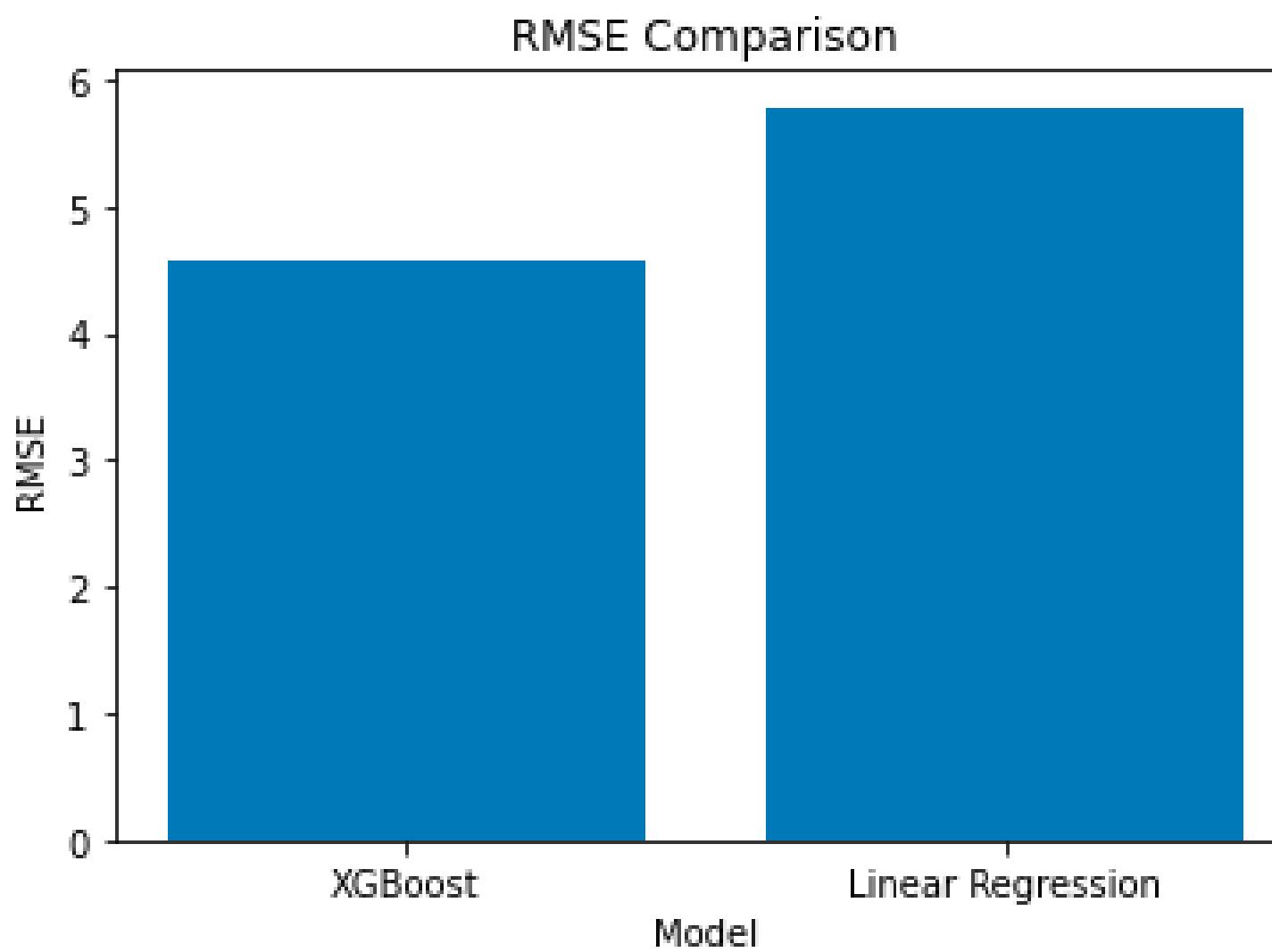
Splitting Training and Testing Data

Fit the model to the data

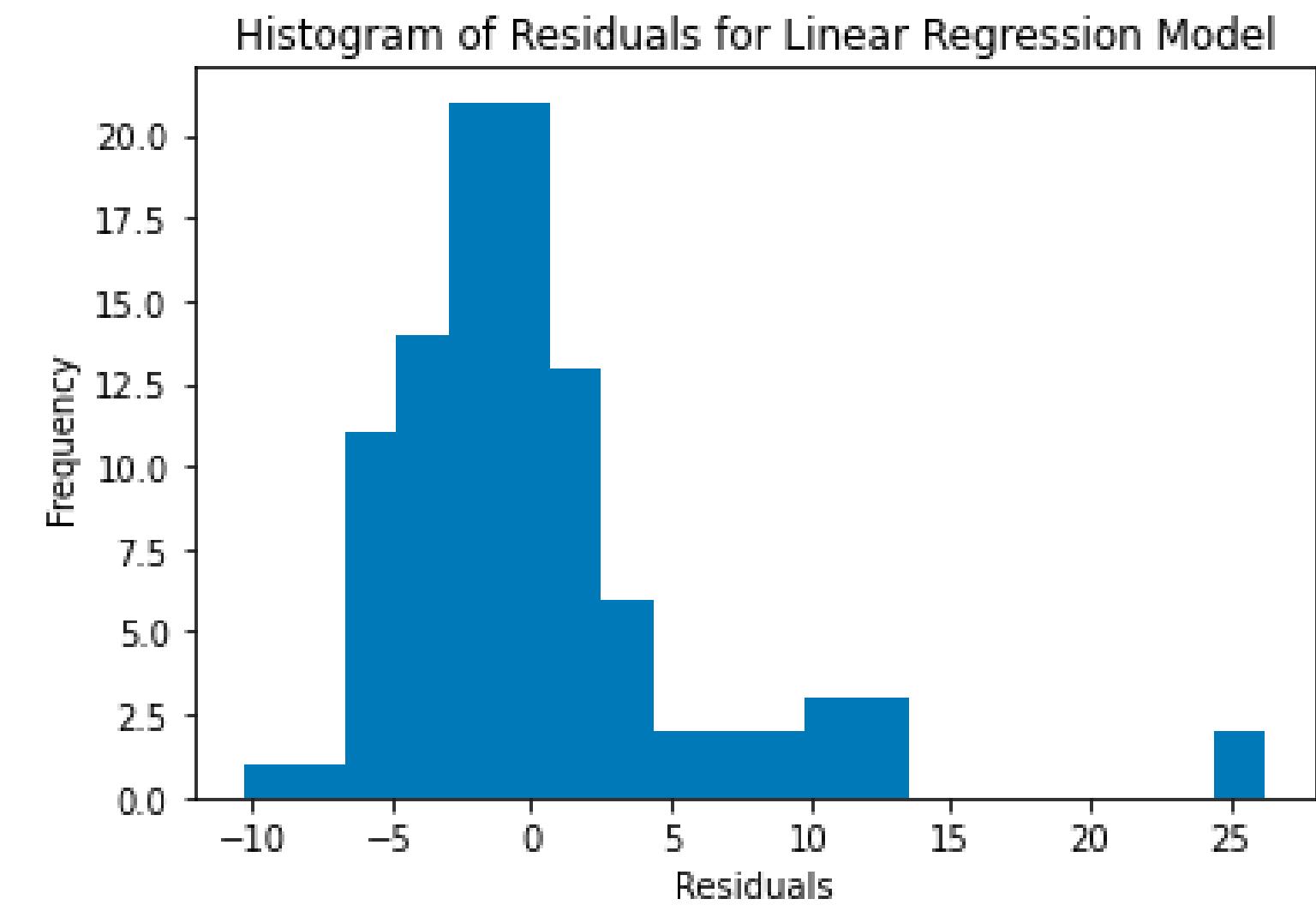
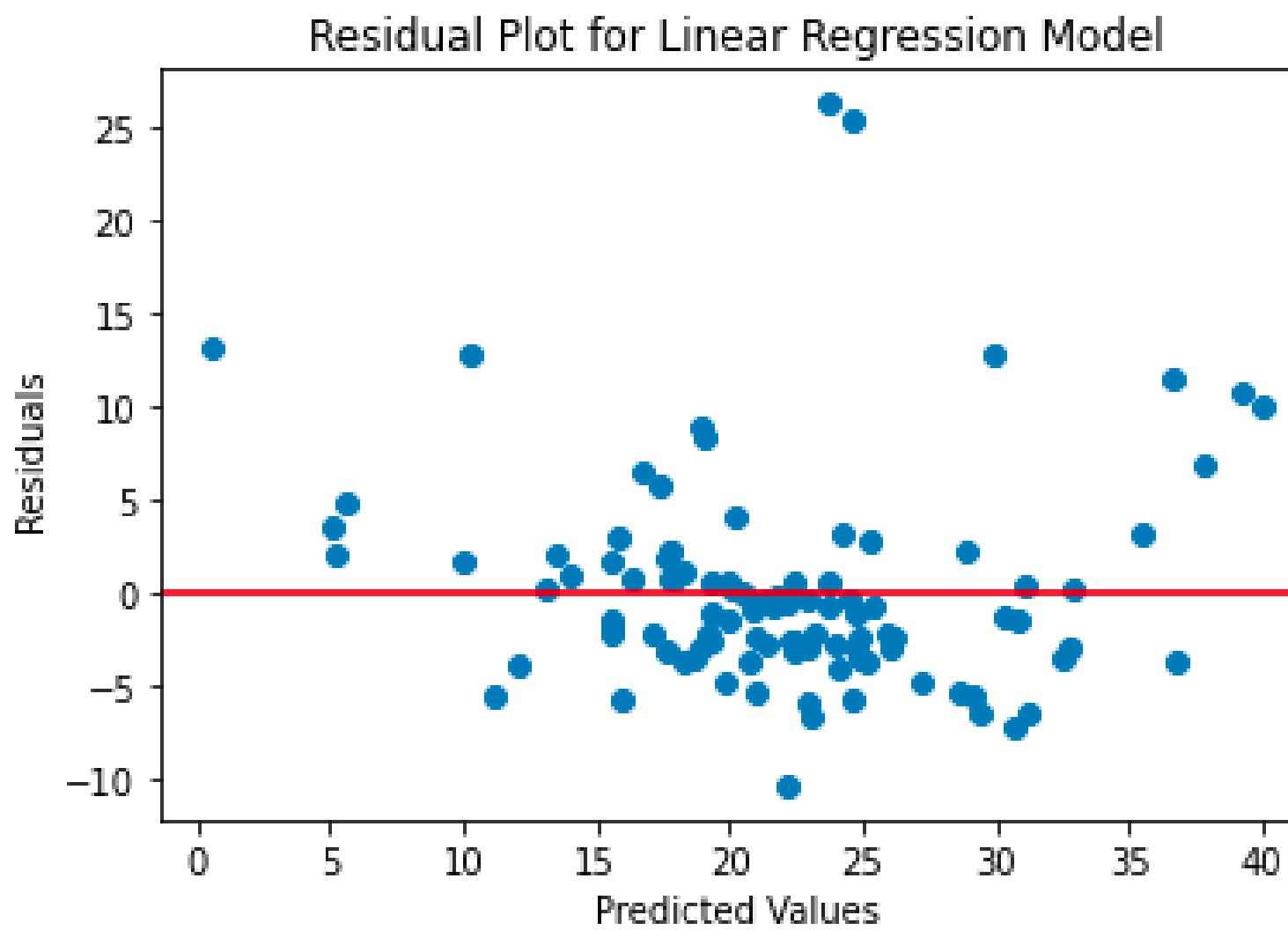
Training the model on the preprocessed data

Deploying the final model to make predictions on new data

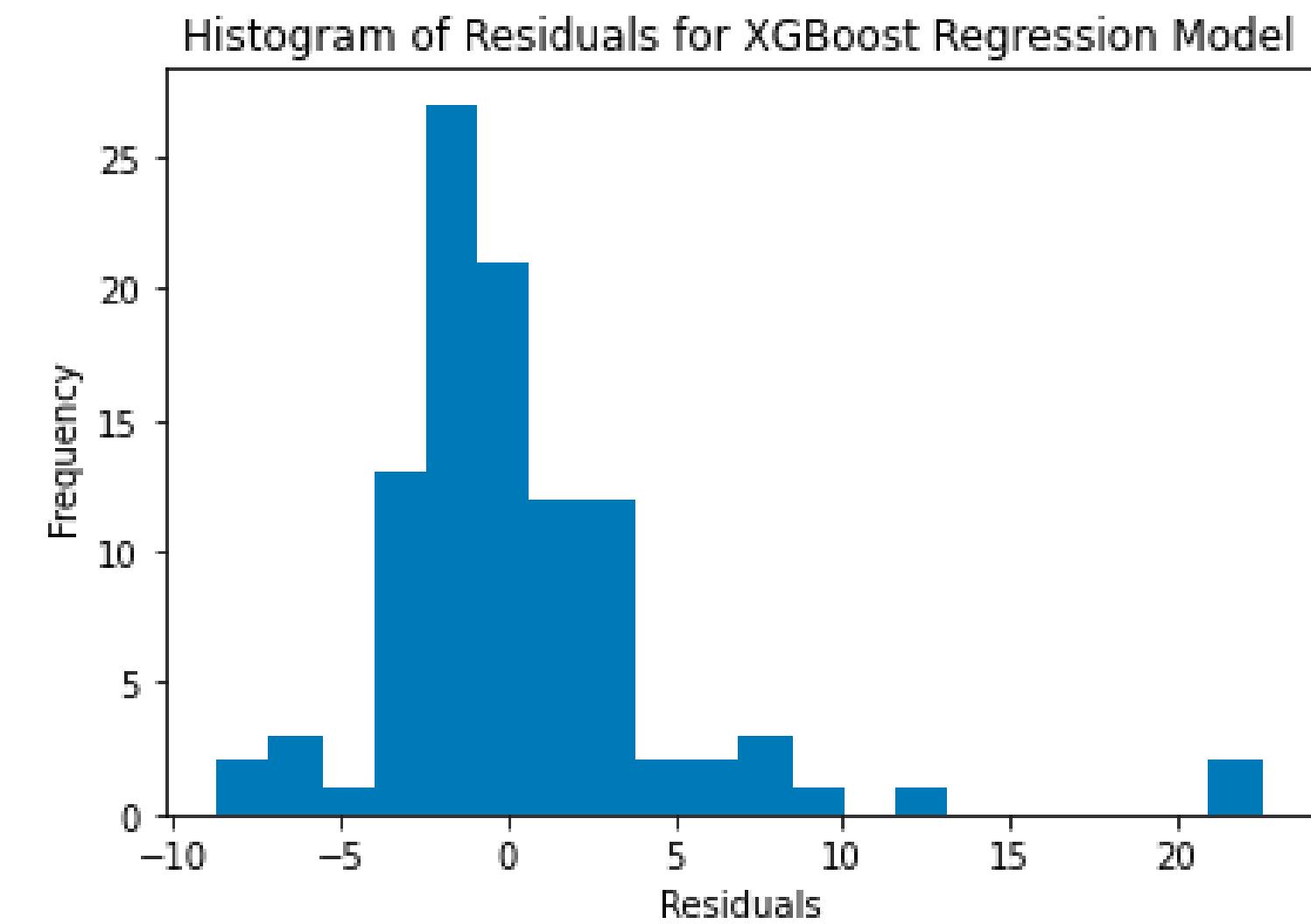
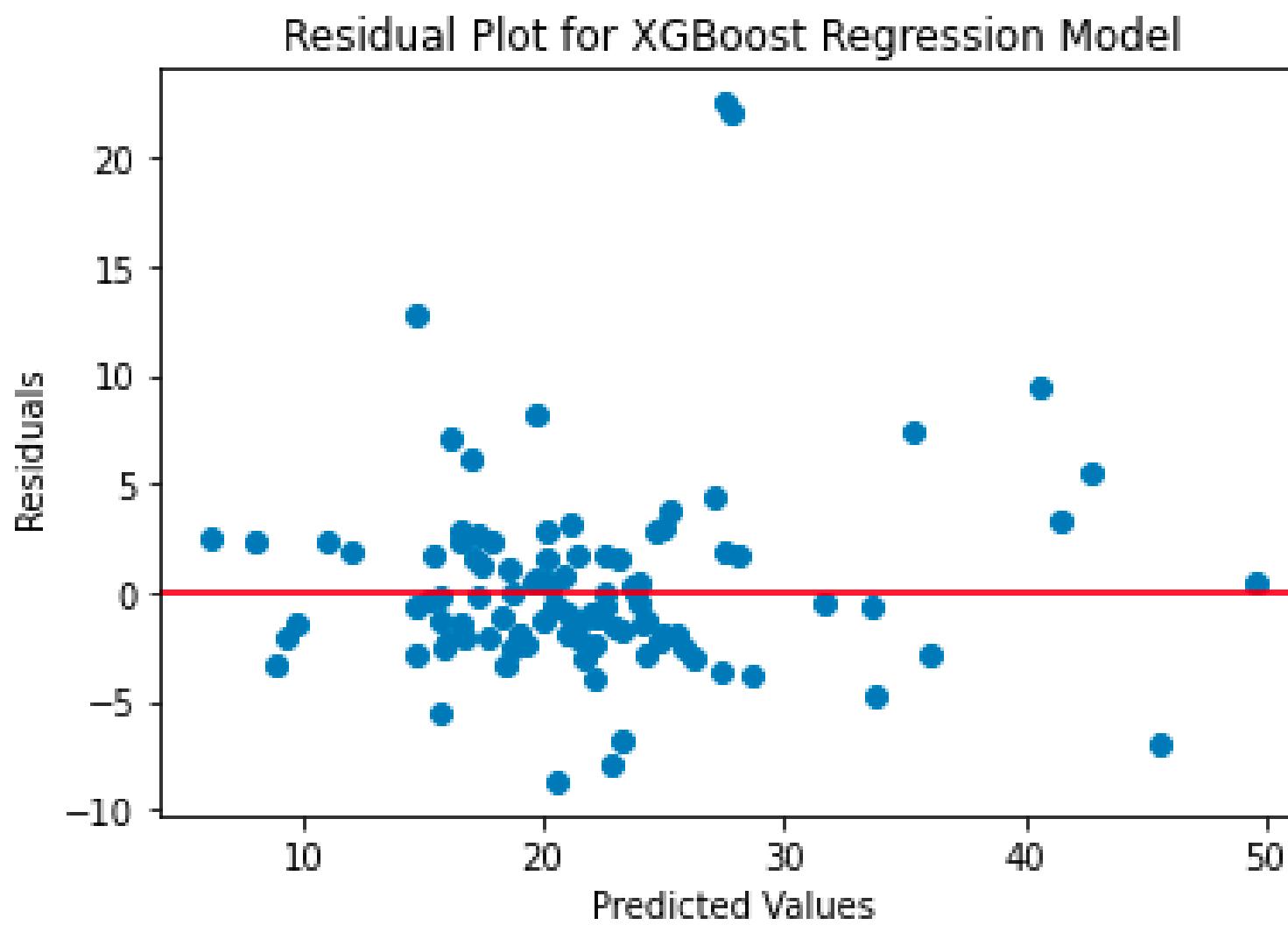
Results



Results



Results



Results

