# HOUSE PRICE PREDICTION USING GRADIENT BOOSTING REGRESSOR, XGBOOST REGRESSOR, AND LIGHTGBM

**UNDERGRADUATE THESIS**

**Submitted as one of the requirements to obtain**

**Sarjana Komputer (S.Kom.)**

**By**
**MUHAMMAD SYAFIQ AYASI**
**001201900054**

**FACULTY OF COMPUTING**
**INFORMATION TECHNOLOGY STUDY PROGRAM**
**CIKARANG**
**SEPTEMBER, 2023**

## PANEL OF EXAMINER APPROVAL

The Panel of Examiners declare that the undergraduate thesis entitled **"HOUSE PRICE PREDICTION USING GRADIENT BOOSTING REGRESSOR, XGB REGGRESSOR AND LIGHTGBM"** that was submitted by **Muhammad Syafiq Ayasi** majoring in **IT** from the Faculty of Computer Science was assessed and approved to have passed the Oral Examination on Thursday September 7, 2023.

**Panel of Examiner**

RUSDIANTO ROESTAM

**Chair of Panel Examiner**

ROSALINA

**Examiner I**

# STATEMENT OF ORIGINALITY

In my capacity as an active student of President University and the author of the thesis/final project/business plan stated below:

Name : Muhammad Syafiq Ayasi

Student ID Number : 01201900054

Study Program : Information Technology

Faculty : Computing

I hereby declare that my thesis/final project/business plan entitled **"HOUSE PRICE PREDICTION USING GRADIENT BOOSTING REGRESSOR, XGBOOST REGRESSOR, AND LIGHTGBM"** is to the best of my knowledge and belief, an original piece of work based on sound academic principles. If there is any plagiarism detected in this thesis/final project/business plan, I am willing to be personally responsible for the consequences of these acts of plagiarism, and will accept the sanctions against these acts in accordance with the rules and policies of President University.

I also declare that this work either in whole or in part, has not been submitted to another university to obtain a degree.

Cikarang, August 21, 2023

(Muhammad Syafiq Ayasi)

## SCIENTIFIC PUBLICATION APPROVAL FOR ACADEMIC INTEREST

As an academic community member of the President's University, I, the undesigned:

Name                 : Muhammad Syafiq Ayasi

Student ID Number  : 01201900054

Study program      : Information Technology

for the purpose of development of science and technology, certify, and approve to give President university a non-exclusive royalty-free right upon my final report with the title:

## HOUSE PRICE PREDICTION USING GRADIENT BOOSTING REGRESSOR, XGBOOST REGRESSOR, AND LIGHTGBM

With this non-exclusive royalty-free right, President University is entitled to converse, to convert, to manage in a database, to maintain, and to publish my final report. There are to be done with the obligation from President University to mention my name as the copyright owner of my final report.

This statement I made in truth.

Cikarang, August 21, 2023

(Muhammad Syafiq Ayasi)

# ADVISOR APPROVAL FOR JOURNAL/INSTITUTION'S REPOSITORY

As an academic community member of the President's University, I, the undesigned:

Name    : Genta Sahuri, S.Kom., M.Kom.

ID number   : 20181000777

Study program: Information System

Faculty    : Computing

Declare that following thesis:

Title of thesis   : House Price Prediction using Gradient Boosting Regressor, XGBoost Regressor, and LightGBM

Thesis of author  : Muhammad Syafiq Ayasi

Student ID Number : 01201900054

will be published in **journal / institution's repository / proceeding / <u>unpublish</u>**.

Cikarang, August 21, 2023

(Genta Sahuri, S.Kom., M.Kom.)

# Final Draft - Muhammad Syafiq Ayasi

or, choose a file to upload

Accepted file types: pdf, docx, doc, txt

☑ I agree to the terms of service                    Get Results

# Your text is likely to be written entirely by a human

## There is a 0% probability this text was entirely written by AI

The nature of AI-generated content is changing constantly. As such, these results should not be used to punish students. While we build more robust models for GPTZero, we recommend that educators take these results as one of many pieces in a holistic assessment of student work. See our **FAQ** for more information.

GPTZero Model Version: 2023-09-14

## Stats

### Average Perplexity Score: 105.024

A document's perplexity is a measurement of the randomness of the text

### Burstiness Score: 122.945

A document's burstiness is a measurement of the variation in perplexity

**Your sentence with the highest perplexity,** *"The Waterfall model is a software develop"*, **has a perplexity of: 671**

# ABSTRACT

The intricate and essential topic of predicting housing prices within the real estate industry will influence customer preferences. Various machine learning algorithms have been applied in an effort to produce predictions that are more accurate. In this work, use and compare the performance of three well-known gradient boosting algorithms for predicting home prices: the Gradient Boosting Regressor (GBR), XGBoost, and LightGBM.

For the training and testing of the model, pertinent housing price information is gathered and compiled. Through rigors feature analysis, significant features for predicting home prices are found and chosen. Additionally, the same dataset was used to train the prediction model using GBR, XGBoost, and LightGBM three gradient boosting methods.

The experimental findings demonstrate that the three algorithms are capable of effectively resolving the issue of house price prediction. Each approach, nevertheless, has benefits and drawbacks in terms of model stability, accuracy, and speed. Based on pertinent evaluation measures, such as R-squared, Mean Absolute Error (MAE), and Mean Squared Error (MSE), we compare the performance of the three algorithms.

In conclusion, picking the appropriate algorithm can enhance the precision and efficacy of house price projections. The findings of this study can serve as a useful manual for practitioners, researchers, and application developers as they select the optimal algorithm for more accurate house price projections.

# ACKNOWLEDGMENT

First, I would like to thank Allah SWT for always giving me the health, strength, and ease that you have given me. Secondly, I would like to thank myself for finally being able to complete my final project. Thirdly, I would also like to thank those who have supported me during my final project support me during the completion of this final project:

1. Mr. Genta Sahuri S.Kom., M.Kom. was the person responsible for supervising my final project, who also helped me and guided me during the final project period.

2. Mr. Rila Mandala Ph.D. as Dean of Faculty Computing, Mrs. Cutifa Safitri, M.Sc, Ph.D. as Program Head of Information Technology, Mrs Indah as Secretary of the Faculty of Computing, and other lecturers who have provided a lot of knowledge during my studies at president university.

3. My beloved parents, for all their prayers and support

4. My friends, who always helped me study during college.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES