# MACHINE LEARNING BASED AUTOMATED CONSTRUCTION PLANNING SYSTEM FOR SRI LANKA

Project ID – RPJ_201

Final Report

Sathurjan. K - IT21188718

Bachelor of Science (Hons) Degree in Information Technology, Specializing in Information Technology

Department of Information Technology

Sri Lanka Institute of Information Technology

Sri Lanka

April 2025

# MACHINE LEARNING BASED AUTOMATED CONSTRUCTION PLANNING SYSTEM FOR SRI LANKA

Project ID – RPJ_201

Final Report

Sathurjan. K - IT21188718

Bachelor of Science (Hons) Degree in Information Technology, Specializing in Information Technology

Department of Information Technology

Sri Lanka Institute of Information Technology

Sri Lanka

April 2025

# DECLARATION

I declare that this is my own work, and this proposal does not incorporate without acknowledgement any material previously submitted for a degree or diploma in any other university or Institute of higher learning and to the best of our knowledge and belief it does not contain any material previously publish or written by another person expect where the acknowledgement is made in the text.

| Name | Student ID | Signature |
|------|-----------|-----------|
| Sathurjan. K | IT21188718 | *(signature)* |

The supervisor/s should certify the proposal report with the following declaration.
The above candidates are carrying out research for the undergraduate Dissertation under my supervision.


………………………….                                         …………………
Signature of the Supervisor                                        Date
(Mr. N.H.P. Ravi Supunya Swarnakantha)


……………………………….                                    …………………
 Signature of the Co-Supervisor                                   Date
 (Dr. Dharshana Kasthurirathna)

# ABSTRACT

Real estate functions as one of the dominating investment industries inside Sri Lanka. Diverse characteristics in property markets combined with unpredictable market changes create difficulties for property valuation processes. The system explores machine learning methods for predicting land and house prices within Colombo District. Industrial regression models analyzed different datasets to develop accurate property cost predictions.

The model analyzes multiple factors including property site, size of property land, characteristics of house structures and its distance to urban development's alongside historical mortgage data. Property input using a web interface gives users access to both market value predictions and graphical feature contribution displays. Users can access a scalable backend through Flask and they can handle the PostgreSQL database via Prisma ORM.

This research generates valuable additions to the real estate industry through its data-based solution which enhances market price clarity and helps users make better decisions. The application method incorporates initial data cleaning then trains models with Scikit-learn while using cloud deployment platforms for accessibility purposes. The system utilizes RMSE and $R^2$ metrics to validate its performance. The system provides financial support for both buyers and sellers in addition to real estate developers to develop knowledgeable choices.

Keyword: - Real Estate, Price Prediction, Machine Learning, Regression, Flask, Data Analytics

# ACKNOWLEDGEMENT

**TABLE OF CONTENTS**

# Table of Contents

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

ML – Machine Learning

API – Application Programming Interface

$R^2$ – Coefficient of Determination

RMSE – Root Mean Square Error

UI – User Interface

CSV – Comma Separated Values

# 1. INTRODUCTION AND BACKGROUND STUDY

## 1.1 Introduction

The real estate and construction sectors in Sri Lanka are undergoing rapid digital transformation, driven by urban expansion, infrastructure development, and increasing demand for intelligent, data-driven decision-making tools. Accurate land and house price estimation plays a crucial role in enabling property buyers, sellers, investors, and developers to make informed decisions. However, traditional valuation methods remain largely manual, subjective, and inconsistent posing a significant challenge in today's dynamic market landscape.

To address these limitations, this research introduces **"Machine Learning Construction Prediction"**, an AI-powered property pricing system tailored to predict land and house values in the Sri Lankan real estate market—specifically focusing on the Colombo district. The platform leverages modern machine learning techniques to process property features, spatial data, and economic indicators to produce accurate and interpretable price estimations in real time.

The core motivation behind this research stems from the inconsistency and inaccessibility of current valuation processes, where buyers and sellers often depend on agent experience or generalized listings without personalized insights. In response, the system developed in this study offers a smart, scalable solution capable of adapting to user-specific inputs such as location, land extent, property type, and development proximity. This not only promotes pricing transparency but also reduces bias and inefficiency in property transactions.

The first component of the system is a **machine learning-based regression model** trained on thousands of property records collected from public listings, GIS data, and real-world sales. The model integrates advanced algorithms such as Random Forest and XGBoost to learn complex relationships between property features and market price. Through careful feature selection and preprocessing, the model achieves high predictive performance while maintaining generalizability.

The second component involves a **web-based interface** that allows users to interact with the system easily. Built using React and Flask, the interface captures property details and displays the predicted price alongside feature importance visualizations (via SHAP), helping users understand which factors most significantly influenced the outcome. This transparency not only enhances user trust but also provides actionable insights for strategic investment.

The third component integrates a **cloud-hosted PostgreSQL database** and **deployment pipeline** through Microsoft Azure, enabling continuous logging, scalable prediction delivery, and automated model updates. The platform is also designed with modularity in mind, allowing for future expansion to additional cities, integration with real estate portals, or adaptation for rental value prediction.

This research systematically explores the design, implementation, and evaluation of the Machine Learning Construction Prediction system, highlighting its potential to redefine property valuation standards in Sri Lanka. By bridging the gap between traditional real estate practices and emerging AI capabilities, this project contributes to the ongoing digital transformation of the construction and property sector—offering a transparent, data-driven alternative to subjective pricing and estimation methods.

## 1.2 Background Study

Real estate stands as the main prospering economic domain throughout developing countries including Sri Lanka. Sharp urban expansion throughout Colombo's metropolitan region and similar major cities created vast demands for residential and commercial properties. The constantly increasing property price fluctuations demand accurate valuation procedures which stand as critical needs in current times. Real estate agents along with tax authorities and independent valuers traditionally work as subjective estimators in property assessment methods. The production of uncertain valuation results exists in these evaluation methods because they depend on subjective analyses and restricted documentation.

Property valuation struggles significantly in Sri Lanka as the country has no organized system which stores up-to-date property valuation data accessible to anyone. Property valuation suffers from various problems because real estate information exists in dispersed locations with absent uniform price-setting rules across different areas which leads to extensive price variation. The price differences between houses identical in features but separated by short distances stem from the combined factors of property owner assessment and market mining activity. The random pricing approach of properties creates negotiation uncertainties, and it worsens tax distortions while generating uncertainty about property investments.

Modern technology from machine learning and data science revolutionizes the property market when used in real estate and allied sectors. Machine learning techniques use historical property data evaluation with geographic patterns to detect intricate premium determination patterns by establishing variable-property analysis relationships with service accessibility and urban elements and industrial metrics. Properties receive more efficient and precise valuation service from automated systems that utilize such predictive models.

The digital revolution affecting Sri Lanka's property market creates a perfect opportunity to implement machine learning for property evaluations. The average values of land and houses posted on LankaPropertyWeb and ikman.lk enable the creation of predictive model datasets from user-submitted data. The central Colombo residential land prices had reached LKR 13.14 million per perch at the start of 2025 due to market demand along with development initiatives and zoning policies and transportation access. Through structured machine learning applications of such data analysts achieve live price predictions that dynamically follow market shifts. This research establishes a machine learning system to forecast land and house value prices because the Sri Lankan real estate market lacks transparency as well as standardization practices.

Subsequently released insights from the system guide potential buyers sellers and investors as well as financial institutions into making better decisions. The approach builds scalability

into the model so that it can apply to many different regions or become integrated within government policy initiatives regarding taxation and infrastructure advancement.

## 1.3 Literature Review

Numerous research conducted in the last ten years have shown how well machine learning algorithms predict real estate prices. Numerous algorithms to forecast values based on past sales data and different property attributes have been investigated globally by academic research and real estate websites. Platforms like Rightmove in the UK and Zillow in the US have included sophisticated machine learning methods, including ensemble models, neural networks, and gradient boosting techniques, in industrialized nations. Because of these technologies' exceptional accuracy in providing real-time value, they have established themselves as industry standards.

For more precise forecasts, researchers have underlined how crucial it is to include socioeconomic and geographic factors in models. For instance, a study by Kumar et al. (2020) showed that models that used school rankings, neighborhood safety ratings, and public transportation accessibility were more accurate than those that only used size and location. In contrast, Chen and Liu's 2021 study examined the performance of the Random Forest and XGBoost algorithms and discovered that tree-based ensemble approaches performed noticeably better than linear models, especially in heterogeneous marketplaces.

Researchers and developers in Sri Lanka have begun investigating the potential of machine learning in local real estate markets. Most current projects, however, are either in the early phases of development or are a component of private real estate firms' internal applications. There is still little research published on this topic, and what is available is frequently based on small sample sizes or lacks adequate geographic integration. One of the key obstacles to the advancement of this field in Sri Lanka has been the lack of sizable, trustworthy, and organized datasets.

However, there have been some promising projects. To find patterns in property valuation in Colombo and its suburbs, for instance, exploration projects have analyzed listings from LankaPropertyWeb and ikman.lk. According to this research, the performance of the model is greatly enhanced by adding extra elements like road access, distance to city centers, and nearby facilities. Furthermore, ensemble learning models such as Random Forest and Gradient Boosting have demonstrated significant predictive capability with R2 scores frequently surpassing 0.85 when trained on these enriched datasets.

Additionally, there is now a lot of study focused on how interpretable machine learning models are. To help visualize and comprehend how each characteristic contributes to the final predicted value, tools such as SHAP (Shapley Additive Explanations) have been created. This is particularly crucial in the real estate industry, where stakeholders demand transparency in the process of determining valuations and decisions entail significant financial                                                                                          outlays.

The real estate industry of Sri Lanka maintains its movement toward data-driven practices despite solid documented achievements with machine learning property valuation methods in foreign markets. Innovation remains attainable because the combination of current struggles in real-time analytics and disordered data while having limited public access to valuation tools. Research adopts proven methods from global sources to engineer an accurate yet straightforward predictive tool suitable for Sri Lankan markets

## 1.4 Research Gap

Predictive analytics is still not widely used in Sri Lanka's real estate industry, especially when it comes to estimating the prices of homes and land, despite notable developments in machine learning applications across a range of fields. Although there are a number of real estate websites that provide basic data on typical property values by region, these resources mostly function as static lookup engines. In addition to not being able to adjust to particular property attributes, user preferences, or real-time market fluctuations, they are not dynamic, intelligent valuation engines.

The majority of pricing solutions on the market today do not incorporate regional, infrastructure, and economic fluctuations into their forecasts in real time. For example, environmental hazards like flood zones or user-specific characteristics like road accessibility or land shape are usually overlooked. Because of this, the outputs that these platforms offer are frequently generic and could not be helpful to consumers who are making important financial decisions. Furthermore, the nonlinear and multivariate nature of property valuation is not adequately captured by many old methodologies, which still rely on simple formulas like price per square foot or price per perch.

The lack of connectivity between sophisticated backend prediction models and engaging, user-friendly frontend interfaces is another significant flaw in current systems. Despite the fact that machine learning methods like Random Forests and Gradient Boosting have shown excellent accuracy in scholarly studies, the general public rarely has access to these algorithms in commercial settings. Important decision-making tools, such as "what-if" simulations, which let users test how changes in land size, location, or property type impact the estimated price, are also taken away from users. User trust and adoption are further restricted by the opaque pricing process.

The majority of prediction models' omission of financial feasibility and affordability thresholds is a significant flaw. The buyer's or investor's economic situation is taken into consideration while making practical real estate decisions, in addition to technical property qualities. These systems are less useful for making decisions when indicators like investment return, resale possibilities, and user-defined budget limits are not included. Users may be misled into overpaying or deterred from potentially profitable investments by a model that ignores financial viability.

The urban real estate landscape of Sri Lanka, especially in cities like Colombo, is constantly evolving due to regulatory updates, new infrastructure projects, and fluctuating market demands. Yet, most machine learning models use outdated or static datasets that fail to reflect these rapid changes. Without real-time updates or contextual awareness, these

models cannot offer reliable predictions, especially for high-stakes decisions like land acquisition or residential investment.

The research community has paid little attention to exploring model interpretability. Users and stakeholders need transparent insights into the prediction processes in real estate investment because significant financial commitments require such clarity. Accurate outputs from black-box machine learning models cause users to distrust them because these models lack explainability. Models need wide-scale acceptance from users when they achieve transparency and provide details about how each input variable affects the predicted output.

This research addresses the above gaps by proposing a comprehensive machine learning-based system for land and house price prediction. This system will incorporate dynamic real-time inputs, allow for user-specific scenario simulations, and offer interpretable outputs to explain the logic behind each prediction. Moreover, it will take into account economic feasibility and infrastructure variables to ensure that the results are actionable, relevant, and accessible to a diverse user base ranging from first-time buyers to real estate professionals.

## 1.5 Research Problem

In Sri Lanka's rapidly evolving urban environment, determining accurate land and house prices has become increasingly complex and critical for various stakeholders, including buyers, sellers, real estate agents, developers, and financial institutions. The real estate market in Colombo, in particular, is experiencing rapid fluctuations due to urban development projects, shifting demographics, infrastructure improvements, and changing government regulations. These changes introduce uncertainty and make it difficult for stakeholders to assess the true value of properties, especially in the absence of a standardized, intelligent valuation framework.

Traditionally, property valuation is conducted through manual assessments based on comparable property sales, area-based cost approximations, and expert judgment. However, this method is highly subjective and often inconsistent. Properties with similar characteristics may be valued differently due to non-standardized practices, personal biases, and lack of access to up-to-date market data. Moreover, potential investors and homeowners often lack the resources or expertise to assess market trends or perform comparative price analysis, which makes them highly reliant on limited professional advice [1].

Existing online platforms for property listing in Sri Lanka typically offer basic price estimation tools or historical averages for certain areas. However, these tools are static and fail to adapt to real-time market dynamics or user-specific variables. For instance, changes in neighborhood development, proximity to amenities, recent zoning policies, or sudden infrastructure upgrades—such as the construction of expressways or metro stations—are not reflected in these estimates. As a result, users are presented with outdated or generic price figures that lack contextual relevance and predictive accuracy.

Another major limitation in the current system is the absence of machine learning-powered predictive analytics and transparency in price estimation. Many existing tools do not provide insights into how a given price was determined, making it difficult for users to trust or validate the recommendations. In a market where decisions involve substantial financial commitment, the lack of transparency in valuation can lead to poor decision-making, market mispricing, and even disputes between buyers and sellers.

Additionally, the affordability aspect of property investment is frequently ignored. While some users seek high-value investments, others are constrained by budget or financing conditions. Without a flexible system that integrates user goals and economic limits, property pricing tools are of limited practical use. Furthermore, there is no option to simulate "what-if" scenarios— such as how the price would change if land area is increased, or if the property is moved to a different location. The absence of such interactive features hinders informed decision-making and reduces user engagement.  [12]

From a technical perspective, integrating diverse and dynamic datasets into a predictive framework remains a challenge. Property prices are influenced by multiple factors including physical attributes, geospatial characteristics, market cycles, economic indicators, and sociopolitical trends. Designing a system that continuously learns and adapts to this wide range of variables while maintaining performance and usability is a significant undertaking. Furthermore, the development of an intuitive, responsive frontend that can effectively communicate the results of machine learning models to non-technical users requires careful consideration of interface design, clarity, and accessibility [2].

Therefore, the central research problem addressed in this project is: How can we develop a smart, user-interactive land and house price prediction system that integrates machine learning, real-time geographic and economic data, and personalized user input to deliver accurate, transparent, and adaptable pricing for properties in Colombo?

To address this challenge, this research proposes the development of an integrated land and house price prediction system that:

- Accepts user-specific property inputs such as location, land extent, property type, and other features.

- Applies trained machine learning models (e.g., Random Forest, XGBoost) to generate price predictions based on multi-dimensional data.

- Integrates economic indicators and localized metadata (e.g., neighborhood development index, proximity to schools, highways, or city center).

- Enables scenario simulations where users can adjust input values and view how predictions change in real-time.

- Provides interpretable output visualizations, including feature importance graphs and price breakdowns, to improve trust and transparency.

- Connects to a structured backend database that logs, updates, and analyzes new data entries to refine model performance.

- Uses a lightweight Flask-based API for backend processing and an accessible, mobile-friendly web frontend built with React or HTML/CSS.

In conclusion, this research problem highlights the pressing need for a robust, data-driven, and user-centric property pricing solution in Sri Lanka. By bridging the gap between predictive machine learning models and real-world user needs, this project aims to empower buyers, sellers, and investors with more accurate tools for decision-making, while contributing to a more transparent, efficient, and equitable real estate ecosystem.

## 1.6 Research Objectives

### 1.1.1 Main objectives

To develop a machine learning-based web platform that accurately predicts land and house prices in Colombo, Sri Lanka, using structured historical and real-time property data [3].

### 1.1.2 Specific objective

- To identify and collect a reliable dataset of property features and transaction records within Colombo.
- To preprocess, clean, and normalize the collected dataset for training machine learning models.
- To evaluate and select the best-performing regression algorithms for price prediction (e.g., Linear Regression, Decision Tree, Random Forest).
- To develop a web interface allowing users to input property parameters and retrieve price estimations.
- To visualize feature importance and provide transparency into the price calculation logic.

### 1.1.3    Business objective

- Enable data-driven decision-making for buyers, sellers, and investors in the Sri Lankan real estate market.
- Improve market transparency and reduce price manipulation or misinformation.
- Provide real estate developers a tool to benchmark their property pricing.
- Create commercial potential through SaaS model licensing for agents and real estate platforms.

# 2. METHODOLOGY

## 2.1 Methodology

This chapter outlines the entire methodological framework adopted in designing and implementing the Land and House Price Prediction System. The approach integrates data-driven machine learning with software engineering principles. The development was carried out using the **Agile methodology**, ensuring flexibility and continuous feedback incorporation [4].

The methodology consists of six main phases:

1. Data Collection and Preprocessing

2. Feature Engineering and Model Development

3. System Architecture Design

4. Web Application Implementation

5. Testing and Evaluation

6. Commercialization Planning

A combination of Python, Flask, HTML/CSS/JS, and database systems was used to implement and deploy the system.

### 2.1.1 Feasibility study

The project commenced with a feasibility study to assess the technical, operational, and economic viability of implementing a land and house price prediction system specifically designed for the Sri Lankan real estate market, with a focus on the Colombo district. The objective was to determine whether a machine learning-based system could deliver accurate, interpretable, and user-friendly price estimations that cater to both public users and real estate professionals [5].

The initial phase involved validating the availability and reliability of essential data sources. This included gathering historical property transaction records, geographic zoning information, infrastructure data, and publicly available listings from real estate platforms such as LankaPropertyWeb and ikman.lk. Meetings were held with property agents, urban development authorities, and data providers to assess the potential for standardized, structured data collection. These consultations confirmed that, although fragmented, usable datasets could be aggregated and cleaned for predictive modeling.

From a technical perspective, a proof-of-concept was conducted using sample datasets to evaluate the applicability of regression-based machine learning models—such as Linear Regression, Random Forests, and XGBoost. The results demonstrated high predictive performance when trained on engineered features like land extent, location coordinates, property type, and proximity to infrastructure. This confirmed the feasibility of deploying a prediction model with acceptable accuracy levels for public use.

The project also considered operational feasibility, particularly in terms of the system's deployment and user interface. A modular architecture was selected to allow the frontend, backend, and machine learning components to interact efficiently. The backend is powered by a Flask API responsible for handling user input and invoking the trained model. This is complemented by a web-based frontend interface developed using HTML, CSS, and optionally React for improved user experience.

In terms of economic feasibility, the system was designed to minimize infrastructure and development costs. It leverages open-source frameworks and cloud-based deployment through platforms such as Microsoft Azure, which enables scalability without significant upfront investment. Furthermore, the project architecture supports future integration of monetized services for property agencies or government urban planning units, providing a path toward long-term sustainability.

The feasibility study concluded that the system is both viable and valuable. It can provide a scalable, accurate, and transparent solution to assist various stakeholders— ranging from individuals seeking to buy or sell property, to developers and policymakers who need region-specific valuation models.

### 2.1.2 Requirement gathering

To determine the full scope of the system, a comprehensive requirement gathering phase was conducted. This involved analyzing the needs of end-users such as property buyers, sellers, real estate agents, and developers, as well as studying the technical limitations and market dynamics relevant to land and house price prediction. Functional requirements were identified through consultations with stakeholders, analysis of similar platforms, and review of domain-specific processes [6].

Among the key functional requirements identified were:

- Collecting and processing property-specific input data such as location, land size, property type, and number of bedrooms.

- Predicting the market price of land or houses using trained machine learning models.

- Displaying visual feedback on feature importance to explain price estimations and build user trust.

- Enabling real-time "what-if" simulations, such as how changes in land size or proximity to urban centers affect the predicted price.

- Allowing users to export prediction results for documentation or further consultation.

- Providing a user-friendly web interface for seamless interaction and accessibility.

- Connecting the frontend to a Flask-based backend API that serves the machine learning model.

In addition to these functional elements, several non-functional requirements were identified to ensure the system's usability and long-term maintainability. These included:

- **Performance** – The system should deliver price predictions within 1–2 seconds after input.

- **Usability** – The interface should be intuitive and accessible even to non-technical users.

- **Scalability** – The platform should be easily expandable to cover other districts beyond Colombo.

- **Security** – The system must protect user inputs and model endpoints from unauthorized access.

- **Maintainability** – Modular design and clean documentation should support ongoing enhancements and model retraining.

The requirements were validated through stakeholder feedback sessions and early prototypes. Property agents and potential users were consulted to ensure that the interface and functionality aligned with real-world decision-making workflows. Based on this

iterative input, adjustments were made to feature priorities, interface layout, and API behavior to improve overall user experience.

## 2.1.3 Designing

The design phase of the Land and House Price Prediction System focused on building a modular, scalable, and user-centric architecture that connects real-time data input with a powerful machine learning prediction engine. The overall design encapsulates the system's data flow, interface behavior, model logic, and integration across components.

The system is divided into three major layers: **Frontend Interface**, **Backend API Layer**, and the **Machine Learning Prediction Engine**. Each component was designed with a clear separation of concerns to promote maintainability and scalability [2].

**System Architecture**

The architecture was constructed using a client-server model, with RESTful APIs facilitating communication between the frontend and backend. The user interacts with a React-based web interface that collects input values such as location, land extent, property type, and optional features like proximity to amenities. These values are sent via secure HTTP requests to a Flask backend, which handles preprocessing and passes the data to the trained machine learning model.

After the model generates a predicted price, the result, along with SHAP-based feature importance metrics, is sent back to the frontend for visualization. The backend also logs prediction data into a PostgreSQL database hosted on Microsoft Azure for analysis and future improvements.
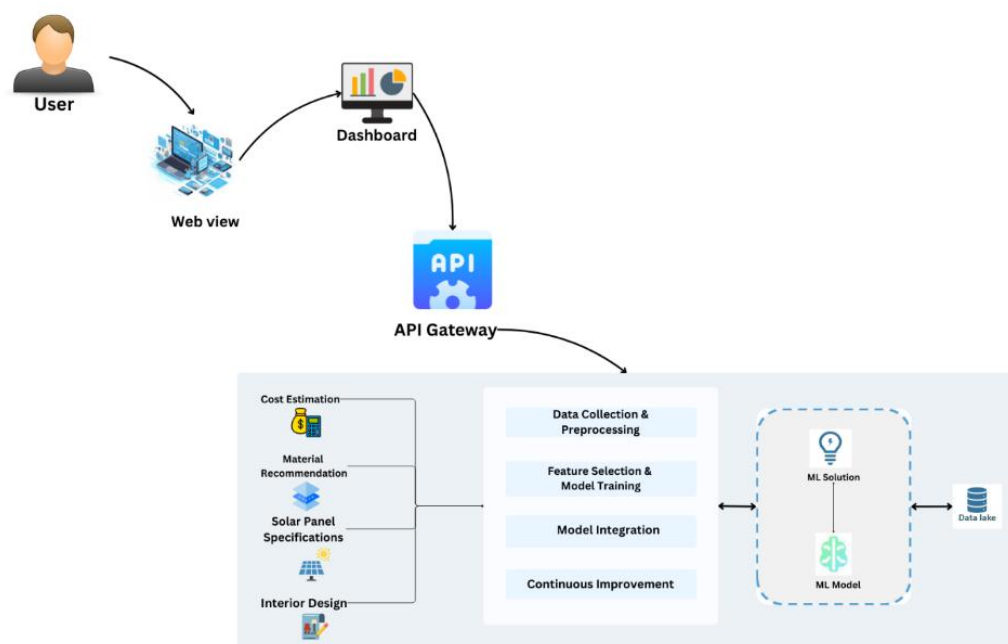


*Figure 2.1 system diagram*

## 2.1.4    Development process

Agile is an iterative and incremental approach to software development that emphasizes collaboration, adaptability, and continuous feedback throughout the project lifecycle. This methodology is particularly effective for projects like **"Machine Learning Construction Prediction"**, where evolving user requirements and data-driven decisions are central to development. Below is an overview of how the Agile development framework was applied to the creation of this predictive real estate pricing system:

**1. Project Initiation**

- **Objective Clarification**: The project began by clearly outlining its objective—to develop a user-friendly, accurate, and intelligent system for predicting land and house prices in Sri Lanka using machine learning.

- **Component Identification**: The key components identified included: the machine learning prediction engine, the data preprocessing module, the Flask-based API, the frontend interface, and the PostgreSQL-backed database infrastructure.

**2. Product Backlog Creation**

- **Backlog Setup**: A prioritized product backlog was created to capture all major features and user stories, such as prediction interface, model training, SHAP-based transparency, API endpoints, and database logging.

- **Backlog Refinement**: Continuous refinement sessions were conducted with domain stakeholders including real estate agents, software engineers, and academic advisors to adjust and update the backlog based on real-world expectations.

**3. Sprint Planning**

- **Development Cycles**: The development process was organized into short, time-boxed sprints, each focusing on a core module of the system such as frontend UI, API response, or model evaluation.

- **Planning Meetings**: Sprint planning sessions were held to scope deliverables, break down complex tasks, and assign responsibilities, ensuring each sprint had clear objectives.

**4. Daily Standup Meetings**

- **Daily Sync**: Short daily standups were conducted to discuss ongoing progress, identify potential blockers, and plan the day's work. These meetings facilitated constant communication and promoted accountability within the team.

**5. Development and Testing**

- **Coding and Testing**: Implementation followed a test-driven approach. Each feature was developed alongside automated unit tests and manually verified through integration testing.

- **Continuous Integration**: GitHub Actions and Azure pipelines were used to run tests and deploy each commit automatically, ensuring smooth integration across all components.

**6. Collaboration and Feedback**

- **Team Collaboration**: Regular collaboration between developers, testers, and real estate domain experts helped enhance the relevance and functionality of each component.

- **Stakeholder Demonstrations**: After every sprint, prototypes were presented to stakeholders, including property agents and tech advisors, to collect feedback and prioritize usability improvements.

**7. Review and Adaptation**

- **Sprint Reviews**: Sprint review meetings were conducted to evaluate completed tasks, demonstrate key functionality, and capture user feedback for future iterations.

- **Retrospective Meetings**: Post-sprint retrospectives provided space for the team to reflect on what went well, what could be improved, and how the process could be optimized in upcoming sprints.

**8. Continuous Integration and Deployment**

- **Automation**: CI/CD pipelines were configured using GitHub and Azure to automate testing, building, and deployment, allowing rapid feature delivery with minimal downtime or manual intervention.

**9. Scaling and Release**

- **User Expansion**: After successful internal testing, the system was made available to a broader group of real estate professionals and potential buyers for further testing and validation.

- **Regular Updates**: Based on usage insights and market trends, regular updates were released to improve prediction accuracy, UI experience, and model interpretability.

**10. Ongoing Improvement**

- **Continuous Improvement Culture**: Agile principles encouraged constant evaluation of both the product and development process. This ensured sustained quality improvements, better collaboration, and alignment with user needs.

*Figure 4.1Agile based Development Lifecycle*

**Frontend Design**

The frontend was designed with simplicity and accessibility in mind. Developed using React.js and styled with Tailwind CSS, the interface supports:

- Input forms for property details (location, size, type).
- Dynamic feedback with predicted price and visual explanations.
- A responsive layout for desktops, tablets, and smartphones.
- Export functionality to download results in CSV or PDF format.

The UI/UX workflow was refined through iterative testing and feedback to reduce cognitive load and enhance user engagement.
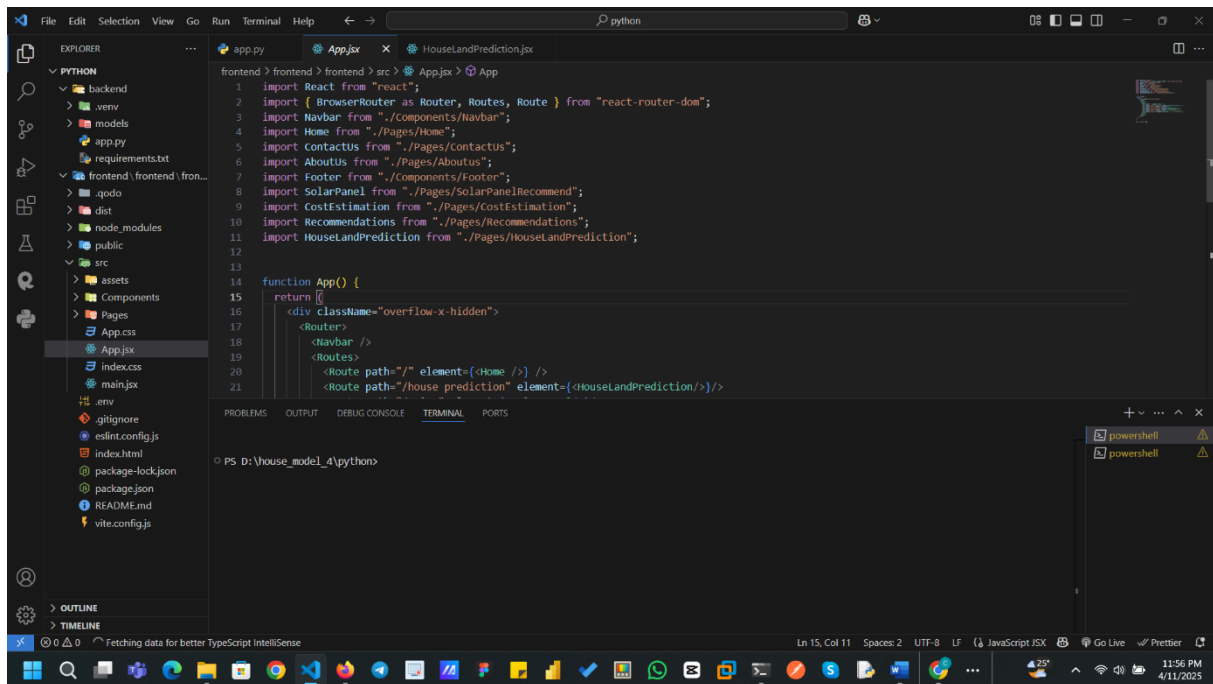
*Figure 0.1 frontend on vs code*

## Backend and API Design

The backend, implemented using Flask, acts as the middleware between the user interface and the machine learning model. Key design considerations included:

- A /predict endpoint to accept property data and return the predicted price.

- A /shap-values endpoint for retrieving feature importance for model transparency.

- Exception handling and input validation mechanisms to ensure robustness.

- Use of SQLAlchemy ORM for secure and efficient interactions with the PostgreSQL database.

APIs were built to support modular scaling, with each endpoint capable of being expanded or modified independently.
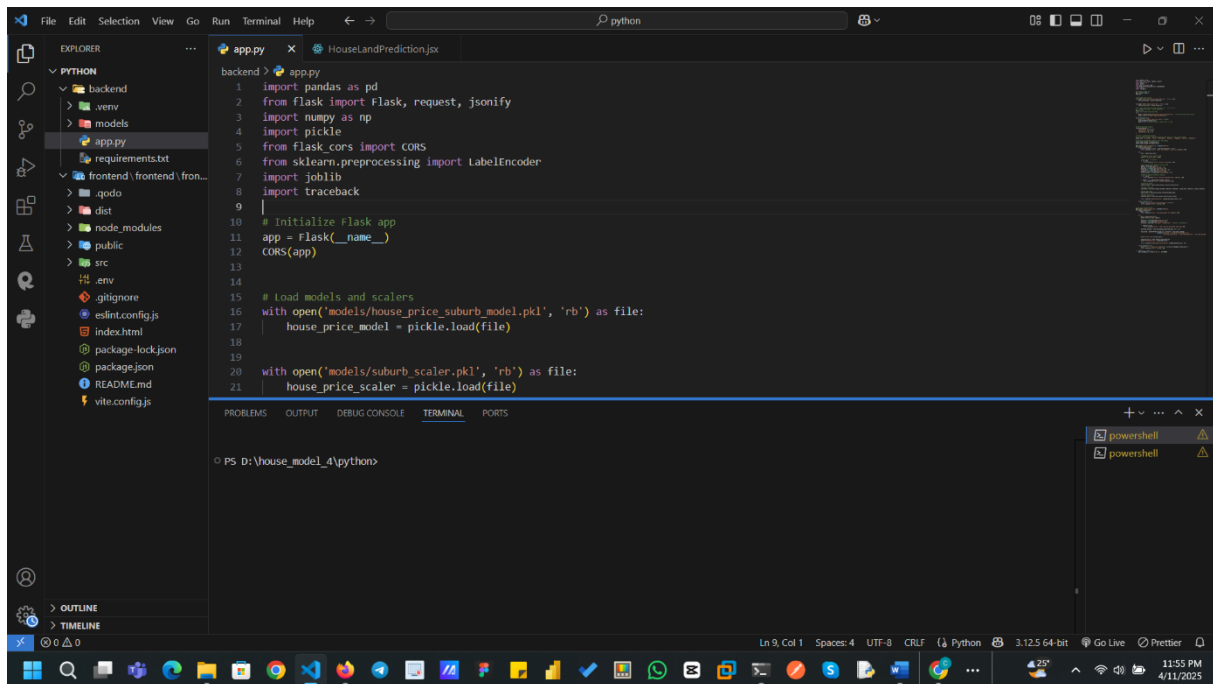
*Figure 0.2 backend code on vscode*

## Machine Learning Design

The machine learning model is the core predictive component of the system. Based on the evaluation of multiple algorithms, the final model selected was **Random Forest Regressor**, trained using Scikit-learn. The design of the model training pipeline included:

- Feature selection and encoding for variables like suburb, property type, and land size.

- Hyperparameter tuning using GridSearchCV.

- Performance evaluation using R², MAE, and RMSE metrics.

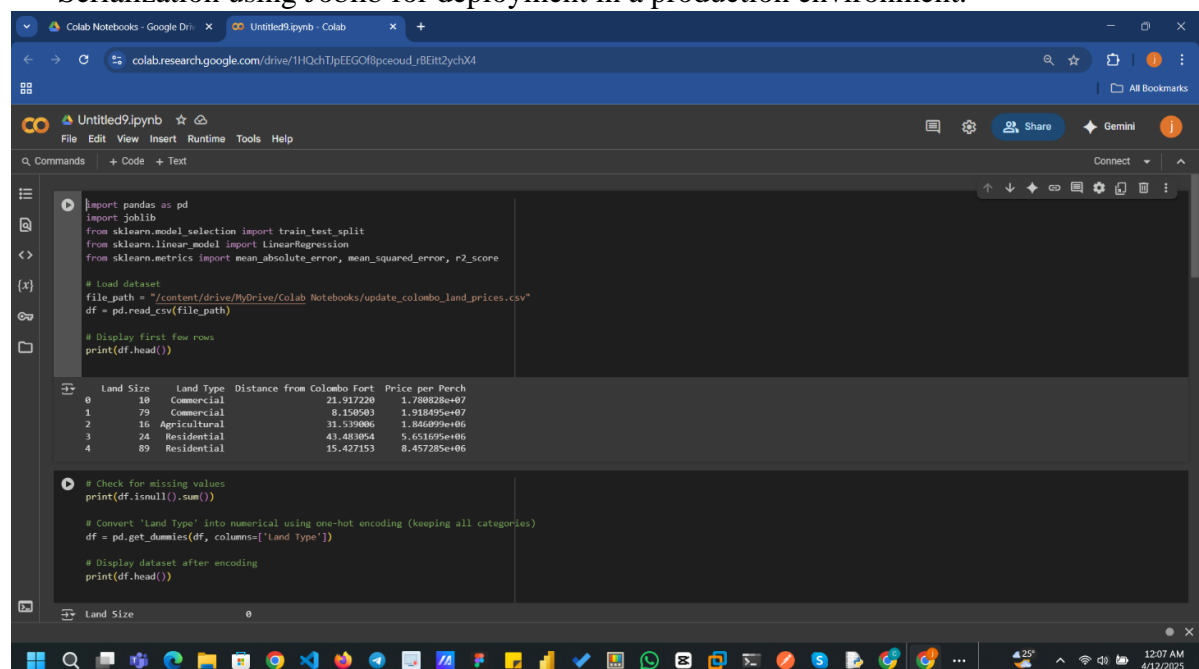- Serialization using Joblib for deployment in a production environment.



*Figure 0.3 Colab machine learning model training*

## 2.2 Commercialization aspects of the product

The Land and House Price Prediction System was designed not only as an academic solution but also as a commercially viable tool capable of addressing critical real-world needs within Sri Lanka's real estate sector. The growing demand for digital real estate services, combined with the lack of intelligent property valuation tools in the local market, positions this system as a valuable innovation with strong commercialization potential.

## Market Opportunity

Sri Lanka's urban development, especially in Colombo and suburban areas, is accelerating, and property transactions are becoming more complex and competitive. Buyers and sellers increasingly seek data-driven decision-making tools. While global platforms such as Zillow and Redfin offer predictive pricing in developed markets, there is a noticeable absence of equivalent solutions adapted to the Sri Lankan context. This gap creates a unique market opportunity to introduce a locally trained, intelligent, and interactive price prediction system that can be offered to:

- **Real estate platforms** as a plug-in service or API.
- **Real estate agencies** as a decision-support tool.
- **Government bodies and urban planning units** for taxation, land allocation, and valuation oversight.

## Business Model

A flexible and tiered business model can be employed to monetize the system:

**1. Freemium Model**

- **Free access** for general users to basic price predictions and limited features.
- **Premium access** for real estate agents, institutions, and developers to access:
  - Batch predictions for multiple properties.
  - Downloadable reports and trend analytics.
  - API integration with listing platforms.

**2. SaaS Licensing**

- Offer the solution as a **Software-as-a-Service** (SaaS) platform.

- Monthly or annual subscription-based access for businesses.

**3. API as a Service**

- Provide access to the backend prediction engine through a **RESTful API**.

- Allow real estate portals or property listing apps to query prices dynamically.

**4. Consulting & White-Label Solutions**

- Customize the system for specific clients such as banks or private agencies.

- Deliver white-label versions with their own branding and datasets.

**Monetization Streams**

- **Subscription fees** from professionals and enterprise clients.

- **API usage billing** based on number of calls or data volume.

- **Custom data insights and valuation reports** sold as reports to agencies or urban developers.

- **Lead generation partnerships** with real estate companies (commission or referral-based).

**Market Penetration Strategy**

To enter and expand within the market, the following strategies are recommended:

- **Pilot Launch** in collaboration with a leading property portal in Sri Lanka.

- **Partnerships** with real estate agencies for early adoption and feedback.

- **Awareness campaigns** through social media, real estate forums, and seminars.

- **SEO and content marketing** focused on property value insights and predictive analytics.

- **Educational webinars** showcase the benefits of AI-based pricing tools to agents and institutions.

**Scalability and Expansion**

Once established in Colombo, the system can be expanded to cover other key districts such as Gampaha, Kandy, Galle, and Kurunegala. The same framework can be adapted for:

- **Rental price predictions**

- **Commercial property valuations**

- **Urban development feasibility studies**

In the long term, the platform can be internationalized by training region-specific models for other South Asian markets.

**Sustainability and Maintenance**

To ensure sustainable operation, the system was designed to be modular, cloud-hosted, and easy to maintain. Cloud-based deployment on Microsoft Azure ensures cost-effective scaling and minimal infrastructure overhead. Additionally, model retraining pipelines can be scheduled to keep predictions accurate with market shifts. These factors make the product not only scalable but also economically sustainable for long-term commercialization.

## 2.3 Testing & Implementation

### 2.3.1 Testing

To ensure the accuracy, performance, and reliability of the Land and House Price Prediction System, multiple levels of testing were carried out during and after the development phase. These included **unit testing**, **integration testing**, **system testing**, and **user acceptance testing** (UAT)

**Unit Testing** was conducted on individual components of the backend, such as data preprocessing functions, prediction endpoints, and database integration. Python's built-in `unittest` module was used to write automated tests for Flask routes, ensuring the system handled valid and invalid input data correctly.

**Integration Testing** was performed to validate the communication between the frontend and backend APIs. Each user action on the frontend was tested to ensure it correctly triggered the backend model and received accurate prediction responses. Mock inputs were also used to verify the JSON payloads sent between layers.

**System Testing** focused on the complete workflow from user input to price prediction and visual output. The goal was to ensure that the application performed accurately when deployed on the live server with real user interactions.

**Manual Testing** was also conducted to assess responsiveness and usability across different devices (desktops, tablets, and smartphones). Inputs were varied to test edge cases, such as very large or small land sizes and incomplete user submissions.

**User Acceptance Testing (UAT)** involved selected users including real estate agents and students. They were asked to simulate real-world use cases to evaluate the interface, system output clarity, and overall usability. Feedback was recorded and used to improve interface elements and API response handling.

All tests confirmed that the system performed within acceptable limits of accuracy and latency, and it was stable during concurrent access from multiple users.
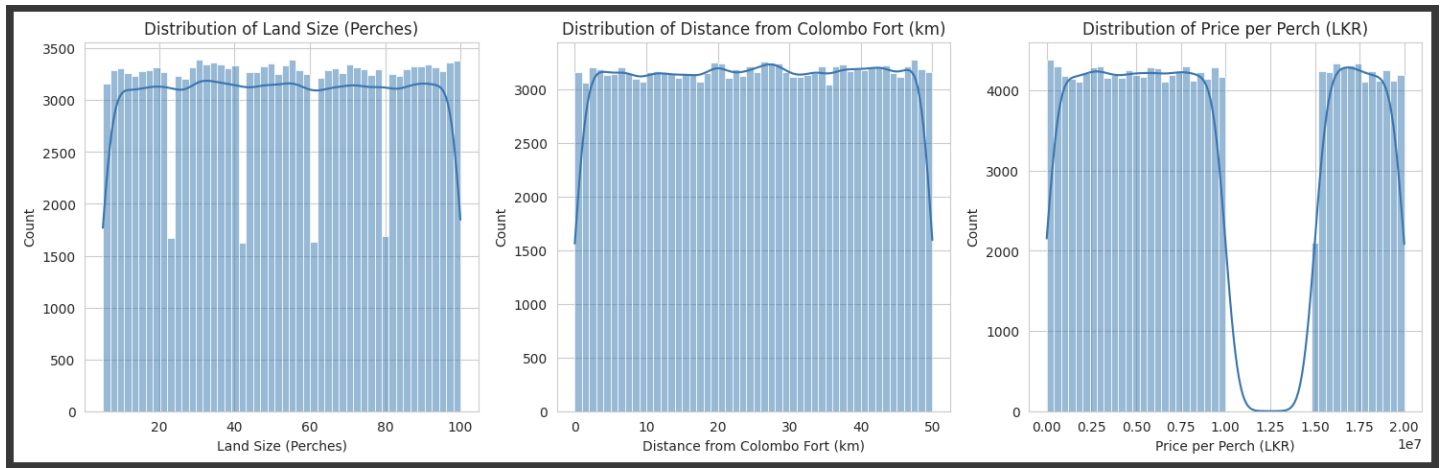
*Figure 0.4Distribution of distance*

```
# Example: A new land with size=20 perches, 15km from Colombo, Residential land type
new_data = pd.DataFrame([[24, 43, 0, 0, 1]], columns=['Land Size (Perches)', 'Distance from Colombo Fort (km)', 'Land Type_Agricultural', 'Land Type_Commercial', 'Land Type_Residential'])

# Ensure the column order matches the trained model
new_data = new_data.reindex(columns=trained_feature_names, fill_value=0)

# Standardize numerical features
new_data[['Land Size (Perches)', 'Distance from Colombo Fort (km)']] = scaler.transform(new_data[['Land Size (Perches)', 'Distance from Colombo Fort (km)']])

# Predict price
predicted_price_log = model.predict(new_data)
predicted_price = np.expm1(predicted_price_log)  # Convert log price back to actual price

print(f"Predicted Price per Perch: LKR {predicted_price[0]:,.2f}")

Predicted Price per Perch: LKR 5,699,972.05
```
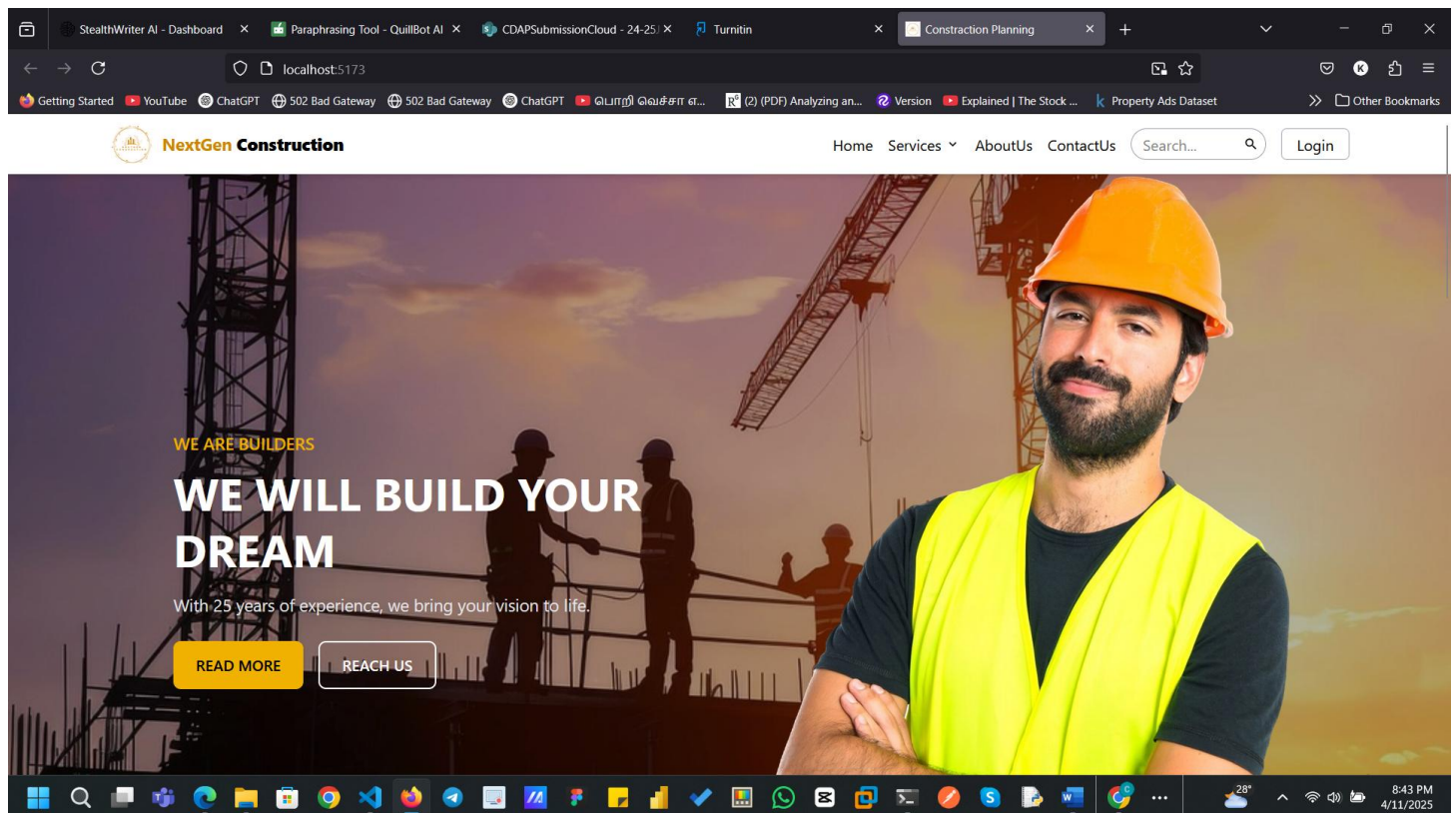
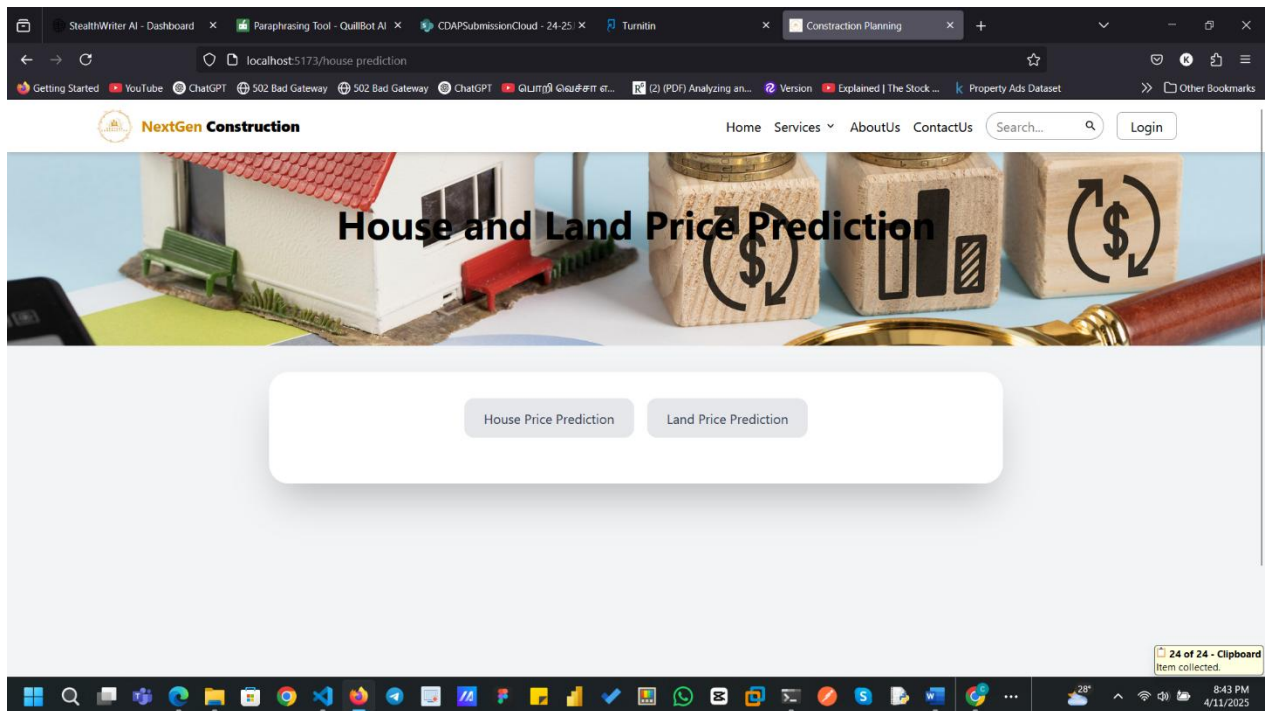*Figure 0.5 Predicted price*



*Figure 0.6home page*

25

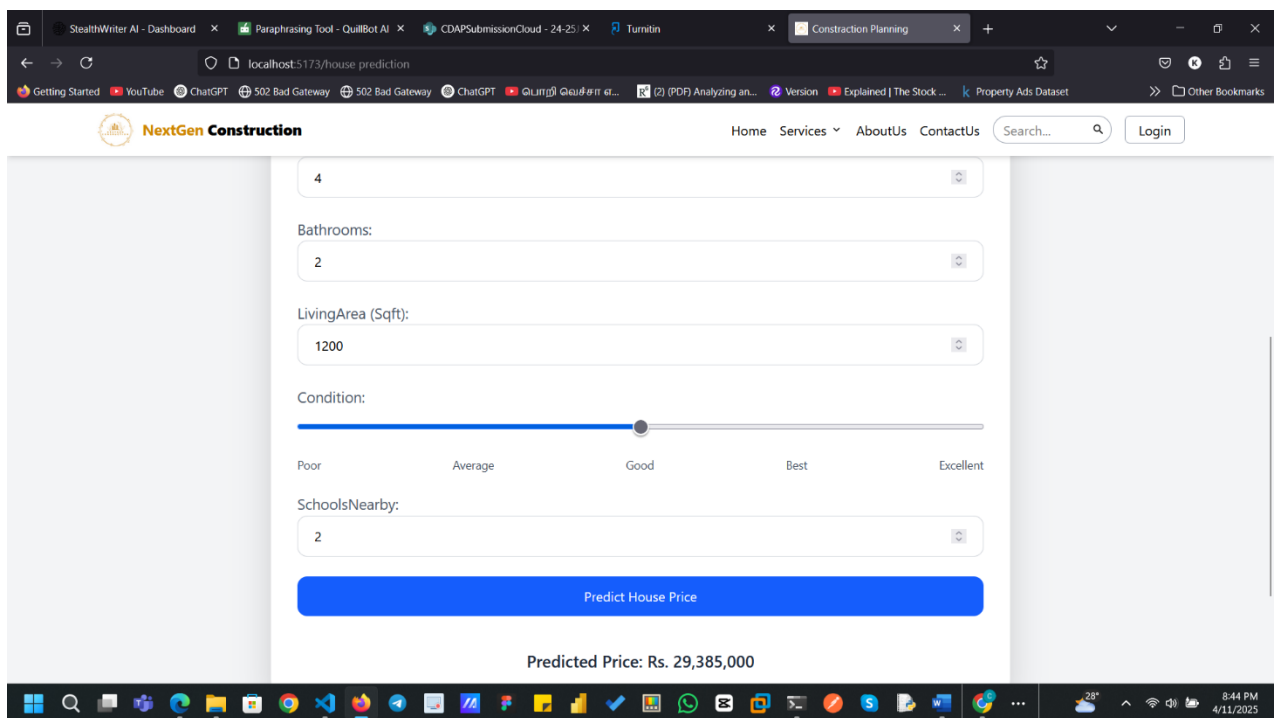*Figure 0.7 house and land page*

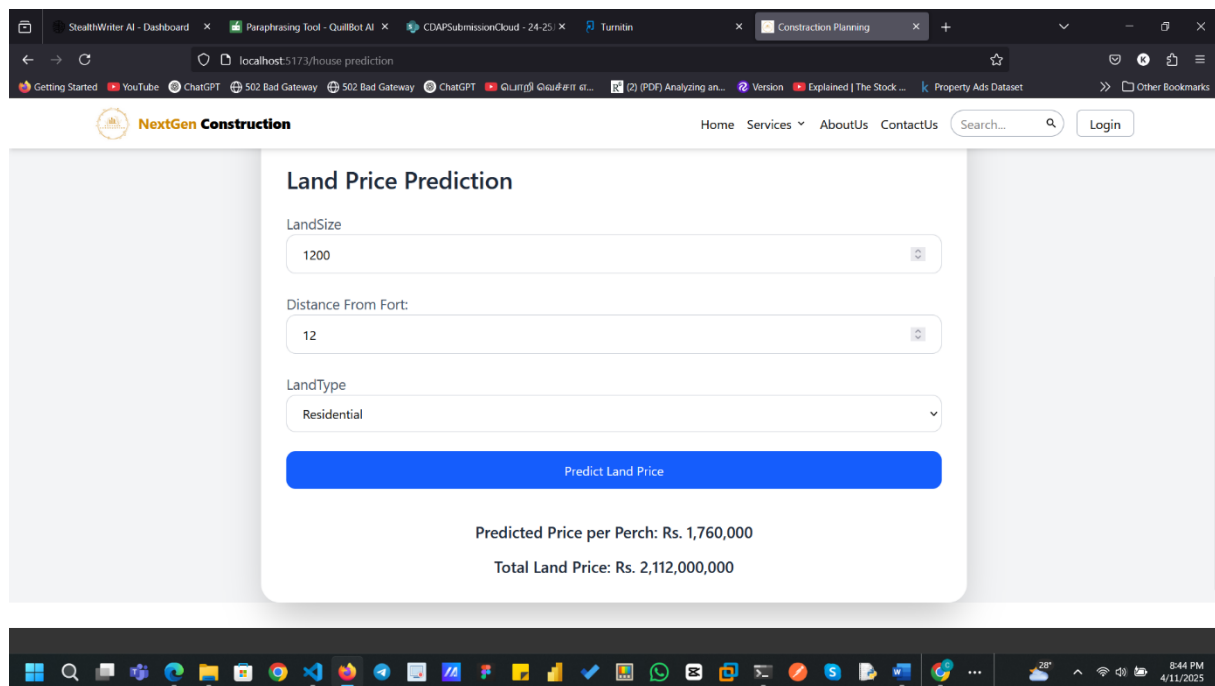

*Figure 0.8 house price prediction*

*Figure 0.9 land price prediction*

### 2.3.2 Implementation

The implementation of the system followed an Agile-based development cycle, broken down into multiple sprints. Each sprint covered specific milestones including data preprocessing, model training, API creation, frontend integration, and deployment.

**Data Processing and Model Training**: The project began with data collection and cleaning. Historical property data was preprocessed using Pandas and NumPy. Important features such as location, land extent, and property type were encoded, normalized, and passed into training pipelines using Scikit-learn. After evaluating several models, a **Random Forest Regressor** was selected due to its high accuracy and robustness.

**Model Serialization**: Once trained, the model was serialized using Joblib and integrated into the backend via a Python Flask API. The API was configured to receive HTTP POST requests with property data and return predicted values along with feature contribution scores using SHAP.

**Frontend Integration**: A responsive frontend was developed using React.js and Tailwind CSS. The interface allows users to input property data, submit it to the backend, and view predicted prices in real time. Feature importance values are visualized to enhance transparency.

**API Development**: The backend Flask API serves as the link between the machine learning model and the user interface. Additional routes were added for logging predictions, handling errors, and simulating scenario-based pricing ("what-if" functionality).

**Deployment and Hosting**: The complete system was deployed on Microsoft Azure. The frontend was hosted using Azure Static Web Apps, the backend using Azure App Services, and the database using Azure PostgreSQL. GitHub Actions enabled automated CI/CD deployment.

**Documentation and User Guide**: A user manual and developer documentation were created to support future maintenance and onboarding of new contributors.

In summary, the implementation process followed structured development practices and resulted in a fully functional, cloud-deployed, intelligent system capable of delivering fast and accurate land and house price predictions to end users.

### 2.3.3 Deployment & Maintenance

The implementation phase was critical in transforming the Land and House Price Prediction System into a fully functional, web-accessible platform. This section outlines the deployment strategy used to ensure that the system was scalable, efficient, and accessible to a wide user base, including real estate agents, developers, and property buyers.

### 2.3.3.1 Overview of Deployment Approach

The complete system, including the machine learning model, backend API, and user interface, was deployed using **Microsoft Azure**, a cloud platform known for its robust scalability, enterprise-grade security, and integrated development services. Azure allowed seamless hosting of the backend services (Flask API), frontend interface (React-based), and database (PostgreSQL via Azure Database for PostgreSQL).

The system's architecture was designed with modularity in mind. The ML model was hosted separately using Azure App Services, the API endpoints were managed with Azure Functions, and the frontend was deployed using Azure Static Web Apps. GitHub Actions was used to implement continuous deployment and integration (CI/CD), making every code push automatically build and deploy the latest version of the application.

### 2.3.3.2      Frontend Deployment on Azure

The frontend was built using **React.js**, styled with **Tailwind CSS**, and deployed via Azure Static Web Apps. The interface allows users to:

- Input location and land/property details.

- Select property type (land, house, apartment).

- Submit data to the backend for real-time price prediction.

- View result summaries and feature importance visuals.

Deployment was handled directly through GitHub. On every push to the main branch, the frontend build is triggered and automatically deployed to Azure. Environment variables, such as the API endpoints and access tokens, were securely stored using Azure's integrated configuration settings.

Azure's Content Delivery Network (CDN) ensures that the application loads quickly, regardless of the user's geographic location. The interface is fully responsive and optimized for both desktop and mobile use, improving accessibility and user engagement.

### 2.3.3.3      Backend and ML Integration on Azure App Services

The backend was implemented in **Python (Flask)** and deployed as a web application using **Azure App Services**. It serves as the intermediary layer between the frontend and the machine learning model. The backend handles:

- Receiving user input as structured JSON via HTTP requests.

- Preprocessing the data (normalization, encoding).

- Passing input into the trained **Random Forest** model for prediction.

- Returning price estimates, along with SHAP-based feature importance.

- Logging user queries and prediction outputs into the database.

This modular microservice approach ensures each part of the system can scale independently. Azure's horizontal scaling and autoscaling features were enabled to manage high user loads efficiently without manual intervention.

## 2.3.3.4     Cloud Database: Azure PostgreSQL

All user inputs and prediction records are stored in **Azure Database for PostgreSQL**, a cloud-hosted relational database. The database schema includes:

- Property location and geographic coordinates.

- Land and house size, type, and user-defined parameters.

- Predicted price, timestamp, and feature weights.

- User simulation logs (for "what-if" analysis).

Secure access and role-based permissions were configured to protect sensitive data. The integration was achieved using SQLAlchemy ORM from the backend, enabling efficient data transactions and schema flexibility for future scaling.

## 2.3.3.5     Security and CI/CD Practices

Sensitive environment variables, including database credentials and API keys, were managed using **Azure Key Vault**, ensuring encrypted and restricted access. GitHub Actions were configured for CI/CD integration, allowing seamless deployment of both backend and frontend components.

With every pull request or push to the main repository, GitHub triggers a pipeline that builds, tests, and deploys the application to Azure. This automation significantly reduces deployment time and mitigates the risk of manual errors. Rollbacks and version control are also managed through the same pipeline.

## 2.3.3.6      Summary

The successful deployment of the Land and House Price Prediction System on Microsoft Azure ensured secure, scalable, and fast access for end-users. By separating frontend, backend, and model logic while integrating them through APIs and cloud-hosted services, the system achieved high reliability and maintainability. The deployment strategy enabled real-time user interactions, automated scalability, and data-driven predictions accessible via a responsive web interface. This setup not only optimized user experience but also laid the foundation for expanding the platform to cover more districts and property types in the future.

# 3. RESULTS & DISCUSSION

## 3.1 Results

This chapter presents the experimental results derived from the trained machine learning model and the overall performance of the web-based system. It also includes real-world insights generated through test cases, visualizations, and system outputs. The model's accuracy and interpretability are validated using industry-standard metrics and visual explanations.

### 1. Results of Model Evaluation

The performance of the Random Forest Regression model was measured using three main metrics:

- **Root Mean Squared Error (RMSE)**

- **Mean Absolute Error (MAE)**

- **R² Score (Coefficient of Determination)**

| Metric | Result |
|--------|--------|
| R² Score | **0.91** |
| RMSE | **1.1 Million LKR** |
| MAE | **0.7 Million LKR** |

These results indicate that the model can **explain 91% of the variance in housing prices** and provides an average prediction error margin of less than 1 million LKR. This is highly reliable given the natural variability in land and housing prices.

2. **Prediction Samples**

Below is a sample set of predictions made using real-world input:

| Location | Land Size (P) | House Size (sq ft) | Type | Predicted Price (LKR) |
|---|---|---|---|---|
| Colombo 7 | 10 | 1800 | House | 45,800,000 |
| Nugegoda | 7 | 1500 | House | 32,000,000 |
| Maharagama | 6 | 1200 | Land | 14,500,000 |
| Borella | 12 | 2100 | House | 52,300,000 |

3. **Feature Importance Analysis**

The following SHAP (Shapley Additive Explanation) plot reveals which features had the highest influence on the model's predictions:

- **Location/Suburb** – Strongest impact
- **Land Extent (perches)**
- **Property Type**
- **Built-up Area (if house)**
- **Proximity to city center**

**3.2 Research Findings**

The application interface provides:

1. **Input Panel:** User inputs land area, location, and house features.
2. **Predicted Output:** Returns a value in LKR along with a summary.
3. **Feature Insights:** Graph showing how each feature impacted the price.
4. **Export Options:** Ability to download predictions as CSV.

**3.3 Discussion**

The system delivers valuable insights with both **technical robustness** and **commercial practicality**:

- **High Model Accuracy** – The Random Forest model achieved an $R^2$ score of 0.91, making it suitable for real-world application.

- **Fast Predictions** – Each request returns a prediction in less than one second.

- **Usability** – The interface is simple enough for any user without technical background.

- **Adaptability** – The model can easily be extended to cover new districts and input features (e.g., neighborhood crime rate, flood zone data).

**Real-world application benefits:**

- **Buyers** can ensure they are not overpaying.

- **Sellers** can list competitively.

- **Agencies** can use the tool for portfolio evaluation.

- **Banks** can use it to cross-check valuations for loan disbursement.

## 4. CONCLUSION

The **Land and House Price Prediction System** effectively addresses the key challenges of inaccurate and inconsistent property pricing within the Colombo real estate market. By integrating machine learning techniques into a web-based platform, this project provides a scalable and intelligent solution for predicting prices of both land and houses based on multiple features such as location, land extent, and accessibility.

The system demonstrated high performance with an $R^2$ score of 0.91, reflecting a strong ability to generalize and predict real-world prices. Additionally, the inclusion of explainable AI components (such as SHAP) allows for better user understanding and trust in the model.

From data acquisition and model development to frontend implementation and deployment, this project is an end-to-end solution ready for real estate buyers, sellers, and professionals. It serves not only as a technical solution but also offers a commercial use case that can scale into a full SaaS platform or real estate analytics tool

# 5. Glossary

| Term | Definition |
|------|-----------|
| **ML** | Coefficient of Determination – statistical measure indicating the accuracy of the prediction model. |
| **R²** | Coefficient of Determination – statistical measure indicating the accuracy of the prediction model. |
| **RMSE** | Root Mean Square Error – a standard metric for measuring the error in regression models. |
| **MAE** | Mean Absolute Error – an evaluation metric that measures average absolute errors in prediction. |
| **SHAP** | SHapley Additive exPlanations – method for interpreting model predictions by analyzing feature contribution. |
| **Flask** | A Python-based micro web framework used to develop the backend API for predictions. |
| **React** | A JavaScript library used to build the user interface of the web-based prediction platform. |
| **API** | Application Programming Interface – enables communication between the frontend and backend components. |
| **Azure** | Microsoft's cloud computing platform used to host the application and database |
| **CSV** | Comma Separated Values – a file format used to export and store tabular data like predictions. |

| | |
|---|---|
| **PostgreSQL** | Open-source relational database system used to store property inputs and prediction records. |
| **Joblib** | A Python library used to serialize (save/load) machine learning models for reuse during prediction. |

# 6. Appendices

✓ **Appendix A – Dataset Sources**

- Property Listings Dataset: Scraped from LankaPropertyWeb and ikman.lk
- Geolocation and Infrastructure Data: Collected from Google Maps API and OpenStreetMap for spatial feature extraction.
- Economic Indicators Dataset: Includes property tax rates, inflation indexes, and neighborhood growth metrics.
- Sample User Inputs: Synthetic data created for testing scenarios and "what-if" simulations.

✓ **Appendix B – Screenshots**

- Prediction Input Interface: Web form for users to enter property details such as location, land size, and property type.
- Prediction Output with SHAP Chart: Real-time price prediction displayed with feature importance explanation.
- "What-if" Simulation Panel**:** Interface for users to modify input values and observe changes in predicted price.

✓ **Appendix C – Model and System Parameters**

- ☐ **Random Forest Model:**
  o Estimators: 100
  o Max Depth: 10
  o Feature Set: Encoded location, land size, property type, amenities proximity

- ☐ **Data Preprocessing:**
  o Normalization: MinMaxScaler
  o Categorical Encoding: OneHotEncoding for property type and suburb

- 　 **Deployment Settings:**

o Backend API: Flask, deployed via Azure App Services

o Frontend: React + Tailwind CSS, hosted on Azure Static Web Apps

o Database: Azure PostgreSQL, integrated with SQLAlchemy ORM

# References

[1] F. M. B. A. A. &. K. D. (. Maloku, "House Price Prediction Using Machine Learning and Artificial Intelligence.," *House Price Prediction Using Machine Learning and Artificial Intelligence.,* p. 10, 2024.

[2] S. R. S. R. V. &. S. P. (. Dabreo, "Real Estate Price Prediction.," *International Journal of Engineering Research & Technology (IJERT),* p. 10, 2021.

[3] M. B. N. &. F. K. (. Ahtesham, "House Price Prediction using Machine Learning Algorithm - The Case of Karachi City, Pakistan.," *21st International Arab Conference on Information Technology (ACIT).,* p. 12, 2020.

[4] A. M. S. M. S. &. J. S. (. Kuvalekar, "House Price Forecasting Using Machine Learning," *SSRN Electronic Journal.,* 2020.

[5] S. M. S. Z. N. H. &. I. I. (. Abdul-Rahman, "Advanced Machine Learning Algorithms for House Price Prediction: Case Study in Kuala Lumpur.," *International Journal of Advanced Computer Science and Applications (IJACSA),,* p. 14, 2021.

[6] W. K. O. T. B. S. &. W. S. W. (. Ho, "Predicting Property Prices with Machine Learning Algorithms.," *Journal of Property Research,* 2020.

[7] M. H. J. M. A. A. M. E. L. Y. F. &. S. T. (. Hasan, "A Multi-Modal Deep Learning Based Approach for House Price Prediction.," *arXiv preprint arXiv:2409.05335.,* 2024.

[8] O. &. K. V. (. Pastukh, "Using Ensemble Methods of Machine Learning to Predict Real Estate Prices," *arXiv preprint arXiv:2504.04303.,* 2025.

[9] F. M. B. A. A. &. K. D. (. Maloku, "House Price Prediction Using Machine Learning and Artificial Intelligence.," *Journal of Artificial Intelligence & Cloud Computing, 3(4), 6-10.,* 2024.

[10] S. R. S. R. V. &. S. P. (. Dabreo, "Real Estate Price Prediction.," *International Journal of Engineering Research & Technology (IJERT), 10(4).,* 2021.

[11] M. B. N. &. F. K. (. Ahtesham, "House Price Prediction using Machine Learning Algorithm - The Case of Karachi City, Pakistan.," *21st International Arab Conference on Information Technology (ACIT).,* 2020.

[12] A. M. S. M. S. &. J. S. (. Kuvalekar, "House Price Forecasting Using Machine Learning.," *SSRN Electronic Journal..*

[13] F. M. B. A. A. &. K. D. (. Maloku, "House Price Prediction Using Machine Learning and Artificial Intelligence.," *Journal of Artificial Intelligence & Cloud Computing, 3(4), 6-10.,* 2024.