



**National Institute of Electronics & Information Technology**

**राष्ट्रीय इलेक्ट्रॉनिकी एवं सूचना प्रौद्योगिकी संस्थान**

**MACHINE LEARNING INTERNSHIP REPORT**  
**HOUSE PRICE PREDICTION**

**Submitted By:**

**Name: Sitasradha Pradhan**

**Program: Foundation Course on Machine Learning  
Using Python**

**Institution: Utkal University**

**Reg. Number: 1782785**

**Submitted To**

**Organisation: NIELIT, BHUBANESWAR**

# Acknowledgment

I take this opportunity to express my profound gratitude to UTKAL UNIVERSITY and the DEPARTMENT OF 5 YEARS OF INTEGRATED MCA for facilitating the successful completion of my internship titled "Foundation Course on Machine Learning using Python."

I am especially thankful to HARIHAR DASH , whose expert guidance, constructive feedback, and consistent encouragement were invaluable throughout the course of this internship. Their mentorship played a critical role in enhancing both my theoretical understanding and practical application of machine learning techniques using Python.

I would also like to extend my sincere appreciation to the coordinators and faculty members associated with the internship programme for designing a structured, industry-relevant learning experience that effectively bridged the gap between academic knowledge and practical implementation.

NAME: Sitasradha Pradhan  
REG. NO: 1782785

# Certificate of Completion

This is to certify that Mr.SITASRADHA PRADHAN,a student of the 5 YEARS OFINTEGRATED M.C.A, UTKAL UNIVERSITY, has successfully completed the internship titled "Foundation Course on Machine Learning using Python", conducted during the period 15<sup>th</sup> May 2025 to 05<sup>th</sup> July 2025.This internship was carried out under the guidance and supervision of HARIHAR DASH, SCIENTIST-C,at NIELIT,BHUBANESWAR.We acknowledge the intern's active participation, consistent effort, and sincere contribution throughout the internship duration. The student has demonstrated a good understanding of machine learning concepts and hands-on proficiency in Python.

**Signature of Guide**

# Table of Contents

	Title	Page No.
	<b>Acknowledgement</b>	
	<b>Certificate</b>	
<b>1</b>	<b>Introduction</b>	<b>1-2</b>
	1.1 Problem Definition and Objective(s) 1.2 Motivation(s) 1.3 Project Overview / Specifications 1.4 Hardware Specification 1.5 Software Specification 1.6 Organisation of the Project	
<b>2</b>	<b>Literature survey</b>	<b>3</b>
	2.1 Existing System 2.2 Proposed System 2.3 Feasibility Study	
<b>3</b>	<b>Python Codes and Project Output</b>	<b>4-8</b>
<b>4</b>	<b>Conclusions</b>	
	<b>References</b>	

# Introduction

## 1.1 Problem Definition and Objective(s) :

House Price Prediction is a regression-based machine learning project that aims to estimate the selling price of a residential property based on various input features such as location, number of rooms, area (square feet), age of the house, and other relevant factors. The system uses historical housing data to train machine learning models that can accurately predict prices for new, unseen inputs. The project is implemented using Python, along with libraries like Scikit-learn, Pandas, and NumPy, and is optionally integrated with Streamlit for a user-friendly interface. The primary objective of this project is:

- Develop a machine learning model that can accurately predict house prices.
- Evaluate and improve the model's performance.
- To predict house prices for new data inputs based on trained models.
- To build and compare multiple machine learning models such as Linear Regression and Random Forest.

## 1.2 Motivation(s) :

In today's fast-growing real estate industry, determining the accurate price of a house is a complex and critical task. Property prices are influenced by a wide range of factors such as location, size, number of rooms, age of the building, and market trends. Relying solely on traditional methods or human estimations often leads to inaccurate pricing, which can result in financial losses for both buyers and sellers.

The motivation behind this project is to create a data-driven solution that uses machine learning to predict house prices more accurately and efficiently. By training models on historical housing data, we aim to eliminate human bias, reduce guesswork, and improve the decision-making process in real estate transactions.

## 1.3 Project Overview / Specifications :

The House Price Prediction project is a machine learning-based system designed to estimate the market price of a residential property based on various input features. The system is implemented using Python and machine learning libraries like Scikit-learn, Pandas, and NumPy.

This project follows a structured machine learning pipeline—starting from data loading and preprocessing to model training, evaluation, and prediction. The primary algorithm used is Linear Regression, but models like Random Forest Regression and Decision Tree Regression may also be implemented for comparison.

The goal is to create a fast, reliable, and accurate system that can assist buyers, sellers, and real estate agents in making informed decisions.

## **1.4 Hardware Specification :**

- Processor: Intel Core i5
- RAM: 8 GB (minimum)
- Storage: 1 GB free disk space
- GPU: Optional for faster training (not required for this project)

## **1.5 Software Specification :**

- Operating System: Windows 10/Linux/macOS
- Python 3.11+
- Required Libraries:
  - pandas
  - numpy
  - matplotlib
  - seaborn
  - scikit-learn

## **1.6 Organisation of the Project :**

- Chapter 1: Introduction
- Chapter 2: Literature Survey
- Chapter 3: Python Codes and Project Output Chapter 4: Conclusions

# Literature Survey

## 2.1 Existing System :

In the existing real estate market, house price estimation is often performed through manual methods, real estate agents, or basic online tools. These traditional systems usually depend on limited data, human experience, or general area-based pricing. As a result, predictions are often inaccurate, inconsistent, and not scalable for large datasets or real-time analysis. Basic statistical methods like simple linear regression have been used but lack the power to capture non-linear relationships and complex patterns in the data.

## 2.2 Proposed System :

The proposed system aims to overcome the limitations of the existing approach by using machine learning algorithms to predict house prices more accurately. The system is built using Python, with the help of libraries such as Pandas, NumPy, and Scikit-learn. Algorithms like Linear Regression, Random Forest Regressions, and Decision Tree Regressions are used to train models on historical housing data.

Once trained, these models can predict house prices based on input features like the number of rooms, location, area (square feet), and age of the property. The system also includes a Streamlit-based web interface that allows users to input property details and receive instant predictions — making it highly interactive, user-friendly, and accessible even to non-technical users.

## 2.3 Feasibility Study :

□ **Technical Feasibility:** The project is technically feasible as it utilizes open-source tools and libraries that are widely supported and documented. It runs efficiently on any modern computer system with basic hardware and does not require high-end infrastructure.

□ **Operational Feasibility:** The system is simple and easy to use, especially with the inclusion of a graphical interface (via Streamlit). Users can enter inputs through form fields and get predictions in real-time, making the system operationally practical and highly usable for real estate agents, buyers, or students.

□ **Economic Feasibility:** This project is economically feasible as all tools used (Python, Scikit-learn, Streamlit, etc.) are free and open-source. There are no licensing costs, and the system can be deployed on local machines or free cloud platforms with minimal expenses.

# Python Codes

```
gui_app.py
import tkinter as tk
from tkinter import ttk, messagebox
import pandas as pd
import joblib

model = joblib.load('house_price_model.pkl') # Change path if needed
root = tk.Tk()
root.title("House Price Predictor")
root.geometry("300x400")
tk.Label(root, text="Choose Location").pack()
location_var = tk.StringVar()
location_dropdown = ttk.Combobox(root, textvariable=location_var)
location_dropdown['values'] = ['Electronic City', 'Whitefield', 'HSR Layout', 'BTM Layout', 'other']
#
Add all locations
location_dropdown.pack()
tk.Label (root, text="Total Sqft").pack()
sqft_var = tk.DoubleVar()
tk.Entry(root, textvariable=sqft_var).pack()
tk.Label(root, text="Bedrooms").pack()
bed_var = tk.IntVar()
tk.Entry(root, textvariable=bed_var).pack()
tk.Label(root, text="Bathrooms").pack()
bath_var = tk.IntVar()
tk.Entry(root, textvariable=bath_var).pack()
tk.Label(root, text="Balconies").pack()
balc_var = tk.IntVar()
tk.Entry(root, textvariable=balc_var).pack()
def predict_price():
    try:
        location = location_var.get()
        total_sqft = sqft_var.get()
        bedroom = bed_var.get()
        bathroom = bath_var.get()
        balcony = balc_var.get()
```



```

if total_sqft <= 0 or bedroom <= 0 or bathroom <= 0:
    messagebox.showerror("Input Error", "Please enter valid non-zero values for Total Sqft,
    Bedrooms, and Bathrooms.")
    return
input_df = pd.DataFrame([
    'location': location,
    'total_sqft': total_sqft,
    'bedroom': bedroom,
    'bathroom': bathroom,
    'balcony': balcony
])
prediction = model.predict(input_df)[0]
price = round(prediction * 1e5, 2)
messagebox.showinfo("Prediction", f"Estimated Price: ₹{price:,.2f}")
except Exception as e:
    messagebox.showerror("Error", f"Something went wrong:\n{e}")
tk.Button(root, text="Predict Price", command=predict_price).pack(pady=10)
root.mainloop()

```

## Housepricepredictor.py

```

import pandas as pd
import numpy as np
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.pipeline import make_pipeline
from sklearn.compose import make_column_transformer
import joblib

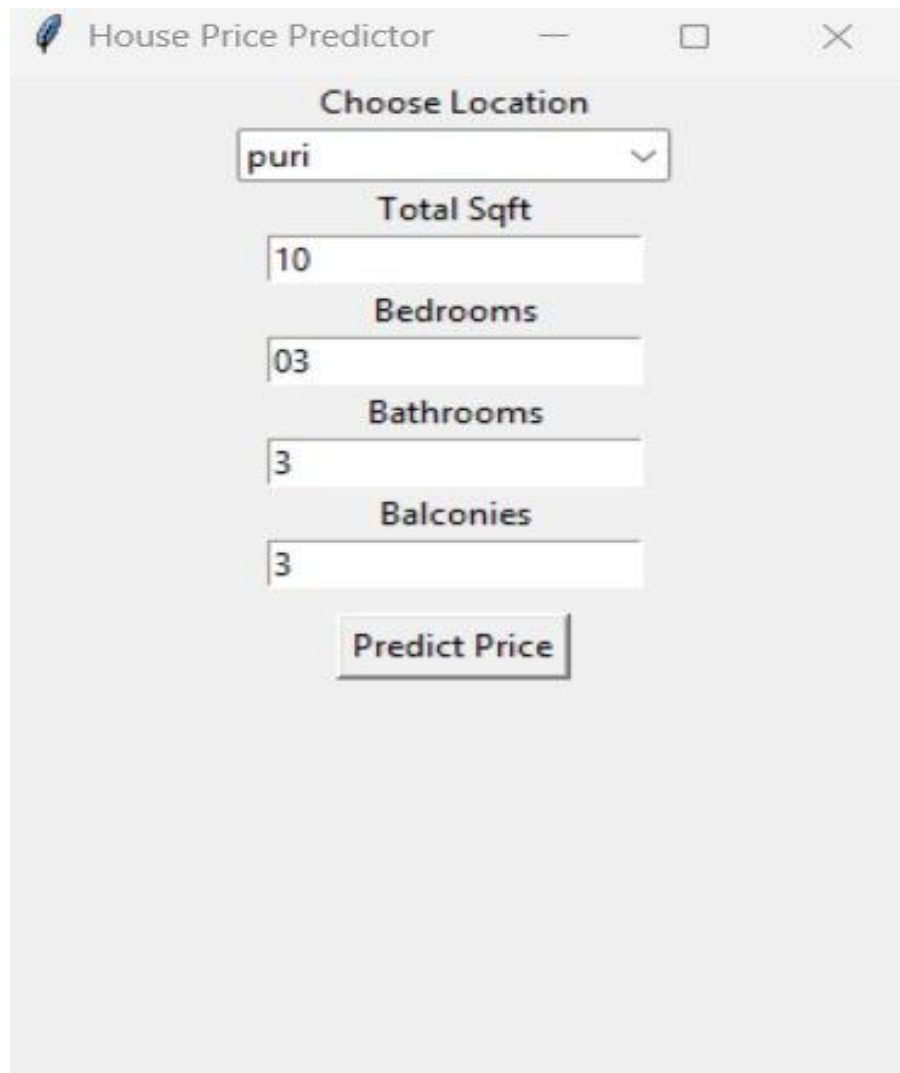
data = pd.read_csv("house_price_large_dataset.csv")
data.drop(columns=['area_type', 'availability', 'society', 'bathroom'], inplace=True)
data.dropna(inplace=True)
data['location'] = data['location'].apply(lambda x: x.strip())
location_stats = data['location'].value_counts()
location_less = location_stats[location_stats <= 10]
data['location'] = data['location'].apply(lambda x: 'other' if x in location_less else x)
data['bedroom'] = data['size'].apply(lambda x: int(x.split(' ')[0]) if isinstance(x, str) else x)
def clean_sqft(sqft):
    try:

```

```
if '-' in str(sqft):
    a, b = map(float, sqft.split('-'))
    return (a + b) / 2
return float(sqft)
except:
    return None
data['total_sqft'] = data['total_sqft'].apply(clean_sqft)
data.dropna(inplace=True)
data['sqft_per_bed'] = data['total_sqft'] / data['bedroom']
data['price_per_sqft'] = data['price'] * 1000 / data['total_sqft']
X = data[['location', 'total_sqft', 'bedroom', 'balcony']]

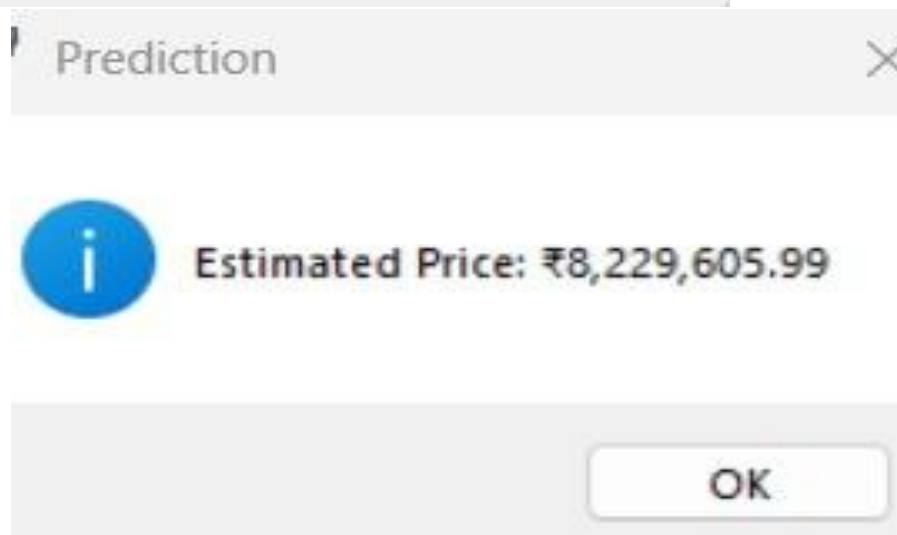
y = data['price']
data.drop(columns=['sqft_per_bed', 'price_per_sqft', 'size'], inplace=True)
col_trans = make_column_transformer(
    (OneHotEncoder(handle_unknown='ignore', sparse_output=False), ['location']),
    remainder='passthrough'
)
model = make_pipeline(col_trans, StandardScaler(), LinearRegression())
x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
model.fit(x_train, y_train)
joblib.dump(model, 'house_price_model.pkl')
data.to_csv('cleaned_data.csv', index=False)
print("Model trained and saved successfully!")
```

# Project Output



A screenshot of a software window titled "House Price Predictor". It contains several input fields and a button. The "Choose Location" dropdown menu is set to "puri". The "Total Sqft" field contains "10", "Bedrooms" contains "03", "Bathrooms" contains "3", and "Balconies" contains "3". A "Predict Price" button is at the bottom.

Field	Value
Choose Location	puri
Total Sqft	10
Bedrooms	03
Bathrooms	3
Balconies	3



A screenshot of a smaller window titled "Prediction". It displays an information icon (a blue circle with a white 'i') followed by the text "Estimated Price: ₹8,229,605.99". An "OK" button is located at the bottom right.

Estimated Price: ₹8,229,605.99

# Conclusions

In this project, we successfully developed a machine learning model to predict house prices using Python. We used essential data preprocessing steps such as handling missing values, encoding categorical variables, and scaling numerical features. The dataset included important features like location, area type, total square footage, number of bathrooms, and availability status.

We trained the model using algorithms such as Linear Regression (or your chosen model like Random Forest, etc.) and evaluated its performance using metrics like Mean Squared Error (MSE) and  $R^2$  Score. The results demonstrated that the model can reasonably predict house prices based on input features. Although we did not use a web interface like Streamlit, the solution works effectively as a standalone Python application, either via command line or GUI (e.g., using Tkinter).

Key Takeaways:

Machine learning can accurately predict real estate prices when trained on clean and relevant data.

Preprocessing plays a crucial role in improving model accuracy. Even without a web-based interface, Python scripts can still be interactive and user-friendly.

# References

- Dataset from Kaggle
- W3Schools. <https://www.w3schools.com>
- Sklearn Documentation. <https://sklearn.io>
- Python Official Documentation , <https://docs.python.org>
- Rajender Kumar. (2020). Python Machine Learning. IK International Publishing House Pvt. Ltd.