# Real Estate Price Prediction Report

Zixiang Huang, Zhongze Zheng, Zhengyu Wang, Yuqing Liu

March 2, 2025

## 1 Introduction

This report presents a detailed analysis and modeling approach for predicting real estate prices using a combination of geographic data, machine learning models, and advanced ensemble techniques.

## 2 Problem Statement

The Toronto real estate market is highly dynamic, influenced by factors such as neighborhood amenities, public transit proximity, economic conditions, and market trends. However, predicting real estate prices accurately remains a significant challenge due to market complexity, missing data, and potential biases. The goal of this project is to develop a machine learning model that can accurately predict real estate prices across Toronto's diverse neighborhoods by analyzing various influencing factors.

## 3 Assumptions for Real Estate Price Prediction

1. **Timeframe of the Data**
   The dataset lacks an explicit collection year, but since real estate prices change gradually, we assume it represents a recent timeframe. Minor misalignment with external datasets (e.g., crime rates) is unlikely to affect model performance.

2. **Crime Data and Trend Consistency**
   Crime rates have shown a consistent yearly increase. Using the most recent crime data (2024) for all listings is reasonable, as past crime levels can be estimated based on proportional growth.

3. **External Data (Infrastructure Locations)**
   Infrastructure data (hospitals, police stations, subway stations) from OSM is assumed to be stable over time, making it a valid representation of public services' availability during the dataset's timeframe.

4. **Model Generalizability**
   The relationship between real estate prices and external factors (e.g., crime rate, infrastructure) is assumed to be stable across similar urban areas, ensuring the model remains applicable for future predictions under similar conditions.

5. **Model Selection (Decision Tree & XGBoost)**

   Housing prices follow a **non-linear** pattern, making Decision Trees more suitable than linear models. XGBoost improves prediction accuracy ($R^2 = 0.91$). We assume the cleaned dataset has no extreme outliers that could hinder model performance.

6. **Economic Factors and Market Stability**

   We assume economic factors (inflation, interest rates) have minimal impact, with housing prices primarily driven by geography and social conditions. However, extreme events (e.g., financial crises) may limit the model's applicability.

# 4 Inputs

| Field Name | Description | Example |
|---|---|---|
| id | Unique six-digit identifier for each listing | 123456 |
| ward | The ward in which the property is located | W10 |
| beds | Number of bedrooms in the property | 3 |
| baths | Number of bathrooms in the property | 2 |
| DEN | Indicates if the property includes a den | Yes |
| size | Property size in square feet | 500-999 sqft |
| parking | Indicates if the property includes a parking space | Yes |
| exposure | Direction the property faces | South |
| D mkt | Days on the market | 30 |
| building age | Age of the building in years | 10 |
| maint | Monthly maintenance fee (C\$) | 500 |
| price | Listing price (C\$) | 750000 |
| lt/lg | Latitude and longitude coordinates | (43.6532, -79.3832) |
| minHighwayDis | Minimum distance to the nearest highway | 1.2 km |
| minMallDis | Minimum distance to the nearest mall | 0.5 km |
| minParkDis | Minimum distance to the nearest park | 0.3 km |
| minPoliceDis | Minimum distance to the nearest police station | 2.1 km |
| minSchoolDis | Minimum distance to the nearest school | 0.8 km |
| minStationDis | Minimum distance to the nearest subway station | 1.0 km |
| hasHospital | Whether there is a hospital nearby | Yes |
| hasMall | Whether there is a mall nearby | Yes |
| hasPark | Whether there is a park nearby | Yes |
| hasPolice | Whether there is a police station nearby | Yes |
| crisisRate | Crisis Rate of Toronto at 2024 | 900 100k/year |

Table 1: Input Features for Real Estate Price Prediction Model

# 5 Mathematical Modeling

## 5.1 Step 1: Initial Data Cleaning and Visualization

We conducted preliminary data cleaning and visualization. Using K-means clustering, we segmented the data into different regions based on geographic coordinates. We observed distinct price patterns in different regions, prompting further exploration of additional influencing factors( Figure 1).
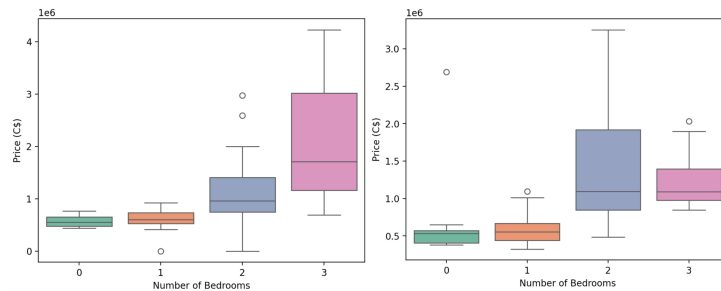


Figure 1: Comparison of the Price Relative to the Number of Bedrooms Based on the Region and Exposure of Houses(Left: region 15, north; Right: region 18, west)

## 5.2 Step 2: Integration of Geographic and External Data

We utilized OpenStreetMap (OSM) to extract geographic data, including infrastructure such as hospitals, police stations, and subway stations (Figure 2). We also integrated crime data from Toronto Open Data, filtering by category (e.g., violent crimes, property crimes) and merging the crime rates into our dataset.
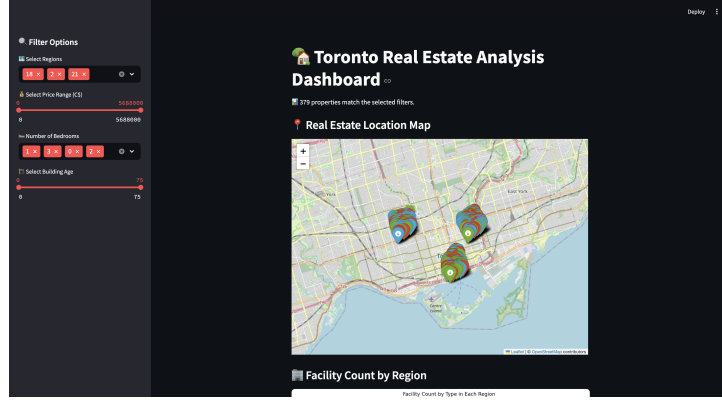


Figure 2: Dynamic Visualization of Housing Information

## 5.3 Step 3: Model Development

We trained and compared various models, including linear regression, lassoCV regression (Figure 3), neural networks, and Random Forest (Figure 4). After comparing accuracy, error metrics, and theoretical alignment, we found that the Decision Tree model best fit our case. We developed a idea of a regional ensemble approach using separate Decision Tree models for each region.
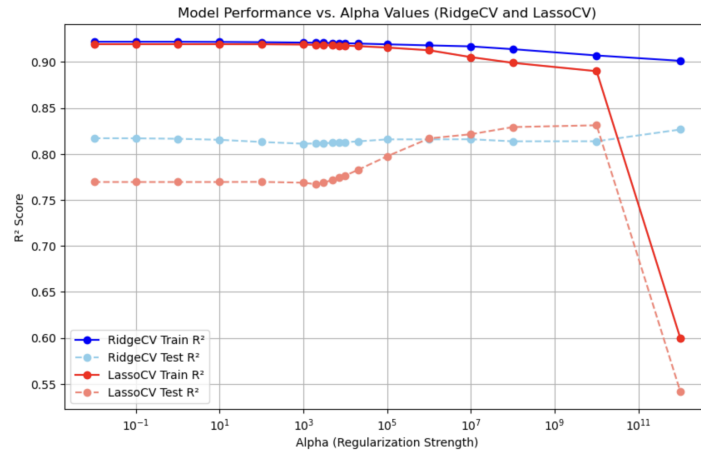


Figure 3: Error of lassoCV and RidgeCV Model with different alpha

## 5.4 Step 4: Model Evaluation and Improvement

We compared a single Decision Tree model with our regional ensemble approach. (Figure 5, 6) The ensemble model outperformed the single model, aligning with XGBoost's principles. Our XGBoost model achieved a promising $R^2$ score of 0.91 on the test set.
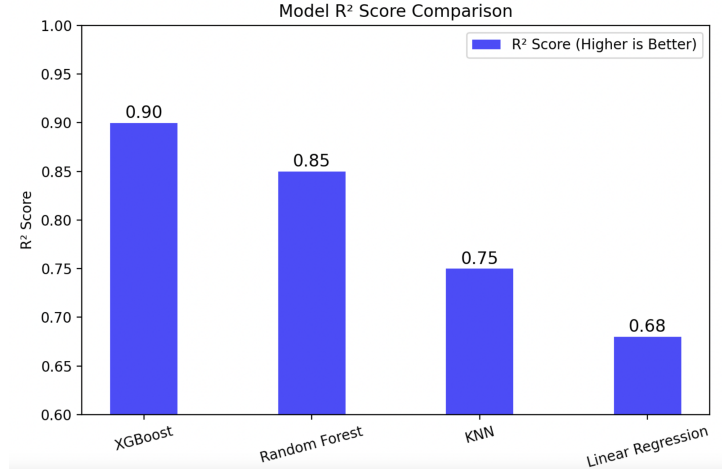
Figure 4: Comparison Between Different Types of Model
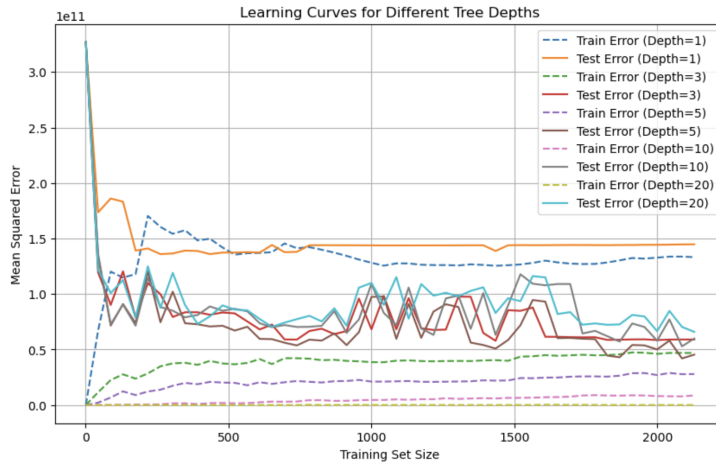


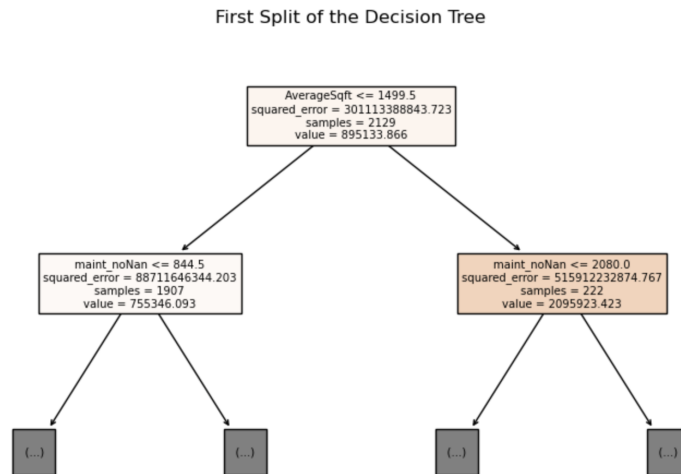Figure 5: Learning Curve of Single Decision Tree With Different Depth



Figure 6: First Split of the Single Decision Tree

# 6 Future and Limitation

From the correlation table between price and all features, our model revealed a feature imbalance. To enhance model diversity and robustness, we propose using the Random Subspace Method for future

```
# Input:
# Training set D = {(x1, y1), (x2, y2), ..., (xm, ym)};
# Base learning algorithm ℒ;
# Number of base learners T;
# Number of subspace attributes d'.

# Procedure:
# for t = 1, 2, ..., T do
#     ℱt = RS(D, d')          # Step 2: Randomly select a feature subset ℱt with d' attributes
#     Dt = Mapℱt(D)           # Step 3: Map the training set D to the feature subset ℱt
#     ht = ℒ(Dt)              # Step 4: Train a base learner ht using the mapped training set Dt
# end for
#
# # Output:
# # The final hypothesis H(x) is determined by majority voting:
# H(x) = argmax (sum_{t=1}^T I(ht(Mapℱt(x)) = y)) for y ∈ 𝒴
```

Figure 7: Implement of Random Subspace Method

adjustments, which could mitigate feature dependency and improve generalization.

# 7   Business Analysis

1. **Data-Driven Property Valuation & Investment Strategy**

   - Real estate investors can use this model to identify underpriced properties in high-potential regions.
   - Homebuyers can compare estimated prices with market listings to avoid overpaying.
   - Developers can predict the impact of infrastructure projects on future property values.

2. **Impact of Crime Rates on Property Prices**

   - City planners & policymakers can use crime data to identify areas needing safety improvements to increase property values.
   - Real estate agencies can integrate crime risk scores into customer reports to help buyers make informed decisions.
   - Investors can find low-crime, undervalued neighborhoods for future appreciation.

3. **The Role of Infrastructure & Public Services**

   - Real estate developers can prioritize projects near high-value infrastructure.
   - Government agencies can plan future infrastructure investments to maximize urban property appreciation.
   - Retail chains & commercial investors can use this data to identify high-value locations for expansion.