

Final Project Part 3 Report

Gadiel David Flores

Wendy Huang

2024-12-05

Contents

Contributions	1
Introduction	1
Methods	2
1. Data Preparation	2
2. Exploratory Data Analysis (EDA)	3
3. Model Building	3
4. Assessment of Assumptions	3
5. Model Diagnostics	3
6. Mitigating Issues	3
7. Conclusion of Methods	3
Results	4
1. Introduction to Results	4
2. Model Comparison	4
3. Key Decisions	4
4. Model Diagnostics and Assumptions	4
5. Interpretation of Final Model	5
6. Visual Representation	7
7. Conclusion of Results	7
Conclusion and Limitations	7
Ethics Discussion	7
References	9

Contributions

- **Gadiel David Flores:** Description of contributions.
- **Wendy Huang:** Description of contributions.

Introduction

Housing affordability has become a pressing issue in many urban cities, including Toronto, where rising house prices have placed both homeowners and tenants out of reach for many years. This paper is conducted to investigating several key house indicators to the influence of house price. This is essential for policymakers , real estate professionals and potential homeowners seeking to navigate this challenging market.

Based on prior experience, we hypothesize that housing supply is a primary determinant of house prices. Therefore, we selected key indicators of housing supply as predictors: the number of available homes, homes sold, homes under construction, new homes completed, and the rate hike, which serves as a proxy for investment conditions. This analysis aims to use the multiple linear regression model to investigate the relationship between house prices and these factors, with the rate hike modeled as a categorical variable capture the impact of interest rate changes on housing affordability.

We conducted literature reviews to validate the selection of our predictors based on prior studies. Housing supply is consistently shown to be a critical factor influencing house prices. Research indicates that higher sales activity (houses sold) is often associated with increased demand, which drives prices upward. Conversely, an oversupply of unsold houses tends to push prices downward. These relationships are commonly examined using hedonic and linear regression models, where transactional data serve as key predictors for estimating price changes effectively (@housepriceprediction).

The number of newly constructed and completed houses influences house prices by affecting market supply. An increase in construction tends to stabilize prices by meeting demand, while lower construction rates can contribute to supply constraints and higher prices. “Global House Prices: Trends and Cycles” highlights the influence of new constructions and completions on house prices. Regression analyses reveal that increases in supply through construction stabilize prices, while supply shortages contribute to price hikes (@housingmarkets)

In the context of rate hikes, the CEPR report discusses their inclusion in econometric models to study the effects of borrowing costs on housing affordability and prices. It finds that higher rates reduce demand, reflected in lower prices, a relationship frequently captured through regression techniques (@whatdriveshouseprices)

Our literature review confirms that the selected predictors—houses sold, houses constructed, houses under construction, unsold houses, and rate hikes—have been previously studied and exhibit a linear relationship with house prices. Linear regression is recognized as a valid method for this analysis, offering a clear way to quantify the relationships between house prices (the response) and these predictors. The validity of the linear relationship will be tested and thoroughly discussed later, as this is a key objective of this paper.

For this analysis, the focus will be on identifying and validating the linear relationship between the predictors and the response variable. We will employ exploratory data analysis (EDA) and evaluate model performance using metrics such as AIC, BIC, and adjusted R^2 . Through these statistical methods, we aim to establish that a linear regression relationship exists between house prices and the predictors: the number of available homes, homes sold, homes under construction, homes completed, and rate hikes.

The structure of this paper is as follows: In the Methods section, we describe the dataset acquisition process and the data cleaning steps undertaken for subsequent analysis. Result section introduce the result of the linear regression and the validation of the model. Limitation section discuss the limitation of the model and dataset selection and processing. Ethics are discussed in the ethics section.

Methods

1. Data Preparation

- Tools Used: Analysis used R with `tidyverse` for data handling and `car` for diagnostics.
- Steps Taken:
 - Defined New Housing Price Index as the response variable and identified predictors like absorption, construction, GDP, and CPI.
 - Removed missing values and duplicates, merged data from 1997-2016, and divided it quarterly.
 - Used Cook’s Distance to remove influential outliers.These steps ensured a clean dataset for accurate regression analysis.

2. Exploratory Data Analysis (EDA)

- Tools Used: `ggplot2` for visualization and `corrplot` for visualizing correlations. Scatterplots were used to check for linearity between predictors and the response variable. Some predictors showed potential relationships and multicollinearity issues. Variance inflation factors (VIFs) flagged these predictors with high multicollinearity for further review. This exploratory data analysis provided essential insights into variable relationships, guiding decisions on transformations and feature selection for improved model accuracy.

3. Model Building

- Tools Used: `MASS` for stepwise regression and `caret` for cross-validation. Both AIC- and BIC-based stepwise regressions tested multiple models to identify the most efficient predictor combinations. Log transformations were used for skewed predictors like `Detached_Unabsorbed_Quarterly_Avg`, improving linearity and stabilizing variance. A Box-Cox transformation ($\lambda = -0.3434$) was applied to the response variable, addressing residual non-normality and heteroscedasticity. These steps balanced model simplicity with predictive accuracy for optimal performance.

4. Assessment of Assumptions

- Tools Used: Residual plots, Q-Q plots, and statistical tests like Cook's Distance. The model's assumptions were thoroughly checked: linearity and homoscedasticity were verified using residuals vs. fitted value and residual vs. predictors plots. Normality of residuals was assessed with Q-Q plots and histograms, revealing minor tail deviations. These steps ensured the regression model's validity, enhancing confidence in its results and inferences.

5. Model Diagnostics

- Tools Used: Performance metrics such as Adjusted R^2 , AIC, BIC, RMSE, and MAE. Model performance was compared using AIC and BIC criteria. The simpler 8-predictor BIC model was chosen for its nearly equivalent performance to the more complex AIC model. These diagnostics ensured the selected model effectively balanced simplicity and performance, adhering to best practices in evaluating regression models.

6. Mitigating Issues

- Tools Used: Log transformations and feature selection techniques. To enhance model reliability, predictors with high VIFs, such as `GDP_Quarterly_Avg`, were excluded to reduce multicollinearity. Log transformations addressed non-linearity in variables like `Completed_Construction_Semi`. Outliers identified through Cook's Distance were removed to minimize their impact on coefficients. These actions ensured the model adhered to regression assumptions and improved its interpretability and accuracy.

7. Conclusion of Methods

The stepwise, structured approach to developing the MLR model ensured robustness and validity while aligning with theoretical principles from the course material. Preprocessing and feature selection techniques addressed potential biases and assumption violations. Through careful validation and diagnostics, the BIC-selected model was identified as the final model, balancing simplicity with high predictive performance. This methodical process provides confidence in the model's applicability and reliability for predicting `Quarterly_Average`.

Results

1. Introduction to Results

This section presents the comprehensive results of the multiple linear regression (MLR) analysis. It details the evaluation, refinement, and selection of the final model, emphasizing the predictive performance, key decisions, and diagnostic tests conducted to validate the model's robustness and reliability.

2. Model Comparison

The analysis involved testing various models to balance predictive power and interpretability. Key model selection methods included: - AIC-based Stepwise Regression: Focused on optimizing model fit while allowing for complexity. - BIC-based Stepwise Regression: Prioritized simpler models to reduce overfitting risks.

Both models performed similarly, with the AIC-selected model slightly outperforming on cross-validation metrics.

3. Key Decisions

To construct a robust and interpretable model, several key decisions were made:

1. Feature Selection:
 - Kept Predictors: Important features like `Detached_Unabsorbed_Quarterly_Avg` and `CPI_Quarterly_Avg` were retained due to their statistical significance ($p < 0.001$) and theoretical importance in explaining price variability.
 - Removed Predictors: Features like `Completed_Construction_Semi` were removed for being non-significant ($p > 0.1$) and to reduce overfitting.
2. Transformations:
 - Log transformations were applied to variables like `Detached_Unabsorbed_Quarterly_Avg` to address non-linearity and stabilize variance.
 - A Box-Cox transformation was applied to the response variable `Quarterly_Average` to improve model fit by addressing skewness and heteroscedasticity.
3. Model Selection:
 - AIC: Used to minimize errors and ensure the model's goodness of fit.
 - BIC: Used to simplify the model and penalize unnecessary complexity, avoiding overfitting.

4. Model Diagnostics and Assumptions

Several diagnostic checks were conducted to ensure the validity and reliability of the model:

1. Residual Analysis:
 - Q-Q Plot: Residuals followed a normal distribution, confirming the assumption of normality.
 - Residuals vs. Fitted Values: Displayed no discernible patterns, confirming linearity and homoscedasticity.
2. Multicollinearity:
 - Variance Inflation Factors (VIFs) flagged high multicollinearity between predictors like `GDP_Quarterly_Avg` and `CPI_Quarterly_Avg`. However, these variables were retained due to their importance in explaining price variability.
3. Influential Points:
 - Outliers and influential points identified using Cook's Distance were removed to avoid skewed or biased results.
4. Validation:
 - Despite the model passing diagnostic checks on training data, test data performance revealed generalization issues, as evidenced by a high MSE and negative R^2 .

5. Interpretation of Final Model

```
library(knitr)

# Model Performance Metrics
performance_metrics <- data.frame(
  Metric = c(
    "Multiple R^2 (Training Data)",
    "Adjusted R^2 (Training Data)",
    "Residual Standard Error (Training Data)",
    "RMSE (Cross-Validation)",
    "MAE (Cross-Validation)",
    "Cross-Validated R^2",
    "MSE (Test Data)",
    "R-squared (Test Data)"
  ),
  Value = c(
    0.9946,
    0.9937,
    0.01029,
    0.01110,
    0.00953,
    0.9933,
    2.621146e+12,
    -8648944050
  )
)

# Coefficients and Predictors
coefficients_table <- data.frame(
  Predictor = c(
    "Detached_Absorption_Quarterly_Avg",
    "Detached_Unabsorbed_Quarterly_Avg",
    "Semi_Unabsorbed_Quarterly_Avg",
    "GDP_Quarterly_Avg",
    "CPI_Quarterly_Avg",
    "Starting_Detached_Construction",
    "Under_Construction_Detached",
    "Under_Construction_Semi",
    "Completed_Construction_Semi"
  ),
  Estimate = c(
    -1.552e-5,
    2.236e-5,
    -2.864e-5,
    -3.004e-7,
    0.01315,
    -4.393e-6,
    6.240e-6,
    2.180e-5,
    -1.250e-5
  ),
  Interpretation = c(
    "Negative effect on housing prices",
```

```

    "Positive effect on housing prices",
    "Negative effect on housing prices",
    "Negative effect on housing prices",
    "Strong positive impact on housing prices",
    "Negative effect on housing prices",
    "Positive effect on housing prices",
    "Positive effect on housing prices",
    "Negative effect on housing prices"
  )
)

# Display Model Performance Metrics Table
kable(
  performance_metrics,
  col.names = c("Metric", "Value"),
  caption = "Model Performance Metrics",
  align = "lc"
)

```

Table 1: Model Performance Metrics

Metric	Value
Multiple R ² (Training Data)	9.946000e-01
Adjusted R ² (Training Data)	9.937000e-01
Residual Standard Error (Training Data)	1.029000e-02
RMSE (Cross-Validation)	1.110000e-02
MAE (Cross-Validation)	9.530000e-03
Cross-Validated R ²	9.933000e-01
MSE (Test Data)	2.621146e+12
R-squared (Test Data)	-8.648944e+09

```

# Display Coefficients Table
kable(
  coefficients_table,
  col.names = c("Predictor", "Estimate", "Interpretation"),
  caption = "Coefficients and Predictors Summary",
  align = "lccc"
)

```

Table 2: Coefficients and Predictors Summary

Predictor	Estimate	Interpretation
Detached_Absorption_Quarterly_Avg	-0.0000155	Negative effect on housing prices
Detached_Unabsorbed_Quarterly_Avg	0.0000224	Positive effect on housing prices
Semi_Unabsorbed_Quarterly_Avg	-0.0000286	Negative effect on housing prices
GDP_Quarterly_Avg	-0.0000003	Negative effect on housing prices
CPI_Quarterly_Avg	0.0131500	Strong positive impact on housing prices
Starting_Detached_Construction	-0.0000044	Negative effect on housing prices
Under_Construction_Detached	0.0000062	Positive effect on housing prices
Under_Construction_Semi	0.0000218	Positive effect on housing prices
Completed_Construction_Semi	-0.0000125	Negative effect on housing prices

The final BIC-selected model included 10 significant predictors. Key metrics:

- Multiple R^2 : 0.9946 (99.46% of the variability explained).
- Adjusted R^2 : 0.9937 (robust against overfitting).
- Residual Standard Error: 0.01029.
- Cross-Validation Results:
- RMSE: 0.01110
- MAE: 0.00953
- R^2 : 0.9933

Significant Predictors:

- Detached_Absorption_Quarterly_Avg ($\beta = -1.552 \times 10^{-5}, p < 0.001$): Negative impact.
- Detached_Unabsorbed_Quarterly_Avg ($\beta = 2.236 \times 10^{-5}, p < 0.001$): Positive impact.
- Semi_Unabsorbed_Quarterly_Avg ($\beta = -2.864 \times 10^{-5}, p < 0.05$): Negative impact.
- GDP_Quarterly_Avg ($\beta = -3.004 \times 10^{-7}, p < 0.001$): Negative impact.
- CPI_Quarterly_Avg ($\beta = 0.01315, p < 0.001$): Strong positive effect.
- Starting_Detached_Construction ($\beta = -4.393 \times 10^{-6}, p < 0.001$): Negative impact.
- Under_Construction_Detached ($\beta = 6.240 \times 10^{-6}, p < 0.001$): Positive impact.
- Under_Construction_Semi ($\beta = 2.180 \times 10^{-5}, p < 0.001$): Positive impact.
- Completed_Construction_Semi ($\beta = -1.250 \times 10^{-5}, p < 0.01$): Negative impact.

6. Visual Representation

- Residuals vs. Fitted Values: No patterns detected, confirming linearity.
- Residuals vs. Predictors:
- Q-Q Plot: Residuals followed a normal distribution.

7. Conclusion of Results

The analysis focused on understanding the factors influencing new home prices, and evaluating a robust model to predict them. Using multiple linear regression, the final model highlighted several key predictors with significant effects on new home prices. For instance, unabsorbed homes had a strong positive impact. Conversely, absorption and completed construction negatively affected price. Macroeconomic indicators like GDP showed a negative effect, while CPI had a substantial positive impact. Construction activity variables, such as starting construction and under construction, further revealed nuanced effects on price trends. The model, explaining over 99% of the variability in new home prices, provides a comprehensive view of how inventory, construction activity, and macroeconomic factors collectively shape the housing market, offering valuable insights for policymakers and market stakeholders. Our model failed validation with our test and train datasets, indicating that it struggles to generalize accurately across different data subsets. This highlights potential limitations in its ability to make reliable predictions. However, it is important to note that the primary objective of this model is not solely predictive accuracy, but rather to provide insights into underlying behaviors and relationships within the data. Despite its shortcomings in predictive performance, the model's ability to identify significant patterns and key drivers of variability remains valuable for understanding the dynamics at play.

Conclusion and Limitations

Ethics Discussion

The dataset used in this analysis originates from Statistics Canada, a reputable and reliable source, which ensures the data's credibility. (StatCan) However, the sensitivity of the response variable—house prices—requires careful ethical consideration. Housing affordability impacts individuals and communities significantly, and any conclusions drawn from this analysis may influence public perception or policy decisions.

First, transparency is essential. We document every step of data cleaning, analysis, and model selection to

ensure reproducibility and avoid misrepresentation of findings. This transparency fosters trust and helps prevent the ethical pitfalls of obscuring methodology. We are fully aware that the results are closely tied to regional economic factors and GDP. We took extra care to avoid errors that could lead to incorrect conclusions.

Second, respect for sensitivity is paramount. House prices are not just economic indicators; they directly affect people’s livelihoods and well-being. It is crucial to present results in a manner that avoids perpetuating inequalities or stigmatizing specific groups or regions. On the other hand, the aim of this paper is to provide insights for policymakers and the estimators. While it is challenging to avoid perpetuate inequalities or stigmatize particular groups or region, we strive to present the results in a way that reflects the data accurately while minimizing harm to any particular group.

Third, fairness in data representation is critical. The predictors—such as sold and unsold houses, construction activity, and rate hikes were analyzed impartially to avoid introducing biases. Any biases identified in the dataset or methodology must be explicitly addressed to maintain the analysis’s integrity.

Finally, as stewards of statistical work, we are responsible for cultivating virtues such as diligence, accountability, and equity, while avoiding vices like negligence or bias. By adhering to these principles, we ensure that the research contributes positively to understanding the impacts of supply-side factors on housing prices in Toronto without causing harm or perpetuating inequities.

References