

Title: Strathcona House Value Predictor

Gilbert Akuja, Tianjiao Jiang, Thamer Aldawood & Yajing Liu

2024-12-07

Table of contents

Summary	1
Introduction	2
Methods & Results	2
1. Import all the necessary libraries for data analysis	2
2. Read in dataset	3
3. Visualization for categorical features	3
4. Prepare data for training and create column transformer for feature trans- formation	4
5. Train, cross validate and evaluate a Ridge regression model	7
6. Predict housing prices with new data	8
7. Visualization for predictions	9
Discussion	9
Limitations	10
Future Work	10
References	12

Summary

Our team will be working on predicting house prices using the **2023 Property Tax Assessment** dataset from Strathcona County Open Data portalCounty (2023). The dataset provides a wealth of information about houses, including attributes like size, location, and other features. By leveraging this data, we aim to build a robust predictive model that accurately estimates house values.

Introduction

The team will be using **Ridge** which is a linear model to predict the value of houses. **Ridge** is a regularization model that is used for predictive modeling and mitigates over fitting, improves model stability especially when features are highly correlated (Hoerl and Kennard (1970)). Ridge helps create robust model that generalize well to new data. The question we aim to answer: Can we predict house prices using publicly available housing data, and which features most influence the predictions? This research is significant as understanding housing prices plays a critical role in assessing broader economic trends, identifying market patterns, and addressing housing affordability issues. Accurate price predictions can help promote transparency and stability in the housing market, which benefits individuals and the economy as a whole.

Data description: For this project we are going to use the 2023 **Property Tax Assessment** from Strathcona County Open Data portal (County (2023)). The data set contains the following attributes related to the different houses. The variables we selected for the model are: **meters** - numeric variable that shows the size of the house **garage** - categorical variable where Y means there is a garage and N means no garage. **firepl** - categorical variable where Y means there is a fireplace and N means no fireplace **bdevl** - categorical variable where Y means the building was evaluated and N means it was not evaluated. The data set was chosen for its rich feature set, adequate sample size, and public availability making it suitable for building a predictive model.

Methods & Results

We used Ridge Regression to predict house values based on features such as building size, garage presence, and building evaluation. Ridge regression, as outlined by Hoerl and Kennard (1970), is particularly useful in addressing multicollinearity. Model selection and evaluation were facilitated by Scikit-Learn's robust tools (Learns (2024)). Exploratory visualizations were created using Altair "Altair Tutorial, Exploratory Data Visualization with Altair" (2024). The report is generated using tools (Quart Team (2024)).

1. Import all the necessary libraries for data analysis

- pandas
- numpy
- sklearn.model_selection
- altair
- os

2. Read in dataset

We read the CSV file named `2023 Property Tax Assessment` into a pandas DataFrame and filters it to include only the features we are going to evaluateCounty (2023). The resulting DataFrame contains the specific features for further analysis. Table 1 displays the first few rows of the dataset after filtering for the relevant columns.

Table 1: Preview of the 2023 Property Tax Assessment dataset after selecting relevant columns.

	meters	garage	firepl	bsmt	bdevl	assess_2022
0	150.59	Y	Y	Y	N	382460
1	123.56	N	Y	N	N	280370
2	104.98	N	N	N	N	402000
3	66.611	N	N	N	N	3690
4	123.83	Y	Y	Y	Y	295910
5	205.4	N	Y	Y	Y	419000
6	120.68	N	Y	Y	N	289380
7	83.148	N	N	Y	Y	326170
8	120.77	N	N	Y	Y	419000
9	115.94	Y	N	N	N	290000

For the beginning exploratory data analysis, we performed data validation checks to ensure the following:

- Correct data file format
- Correct column names
- No empty observations
- Missingness not beyond expected threshold
- Correct data types in each column
- No duplicate observations
- No outlier or anomalous values

The results for our validation check are passed.

3. Visualization for categorical features

We have enabled VegaFusion to optimize Altair data transformations for visualizations“Altair Tutorial, Exploratory Data Visualization with Altair” (2024). We then created a distribution plot of categorical features in the `housing_df` DataFrame. As shown in Figure 1, the distribution of categorical features provides insights into their overall counts.

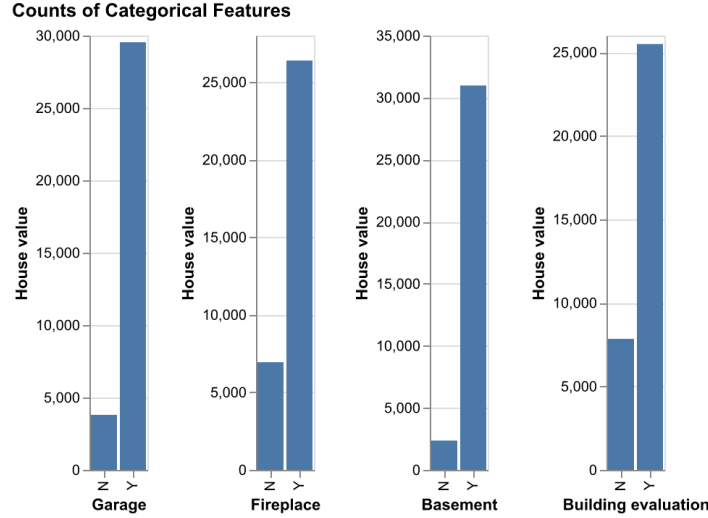


Figure 1: Distribution of categorical features (Garage, Fireplace, Basement, and Building Evaluation).

The scatter plots visualize the relationship between house assessment values (`assess_2022`) and four categorical features: `garage`, `firepl`, `bsmt`, and `bdevl`. As depicted in Figure 2, the scatter plots illustrate the relationship between house value assessments and categorical features.

The correlation plots show which of our features have positive or negative relationships with each other. By examining Figure 3, we can see that our features have no negative relationships, and that there seems to be a strong positive relationship between `meters` and `assess_2022` which tells us that the larger the property is, the higher its evaluation. We can also see a significant positive relationship between `bsmt` and `bdevl` which shows that most properties that have a basement in our dataset have had a building evaluation.

4. Prepare data for training and create column transformer for feature transformation

Visualizing the distribution of the target variable, `assess_2022`, using a histogram and Kernel Density Estimate (KDE) plot is a crucial step in data exploration. It helps us understand the overall distribution, identify patterns such as skewness, kurtosis, or multiple peaks, and detect outliers. By comparing the histogram and KDE, we can assess if the distribution meets the assumptions of our chosen modeling techniques. This visualization also guides decisions about preprocessing, such as applying transformations to address skewness or choosing the right algorithm for modeling. Additionally, it can uncover data quality issues, like sparse regions or missing values, prompting necessary corrections.

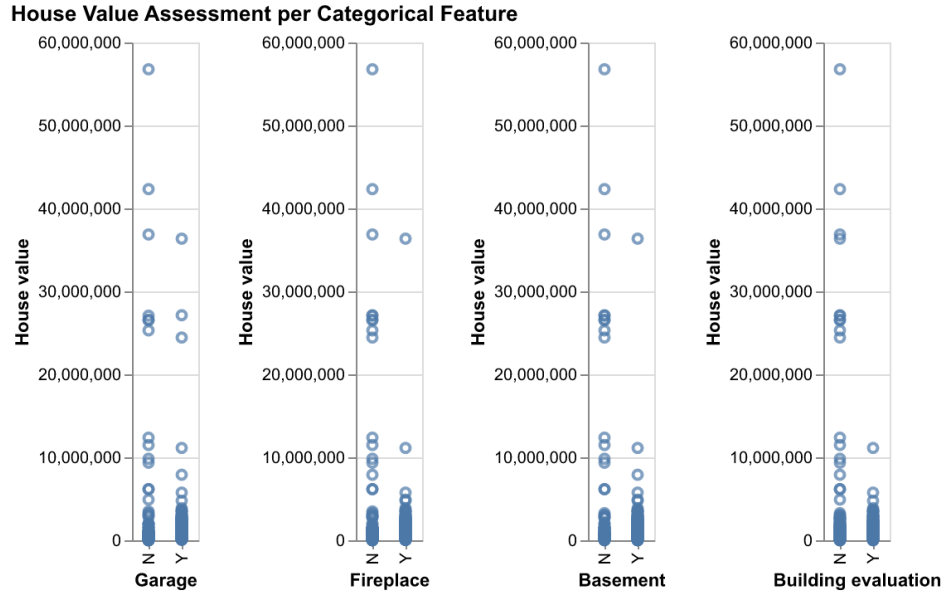


Figure 2: Scatter plots showing house value assessments for categorical features: Garage, Fireplace, Basement, and Building Evaluation.

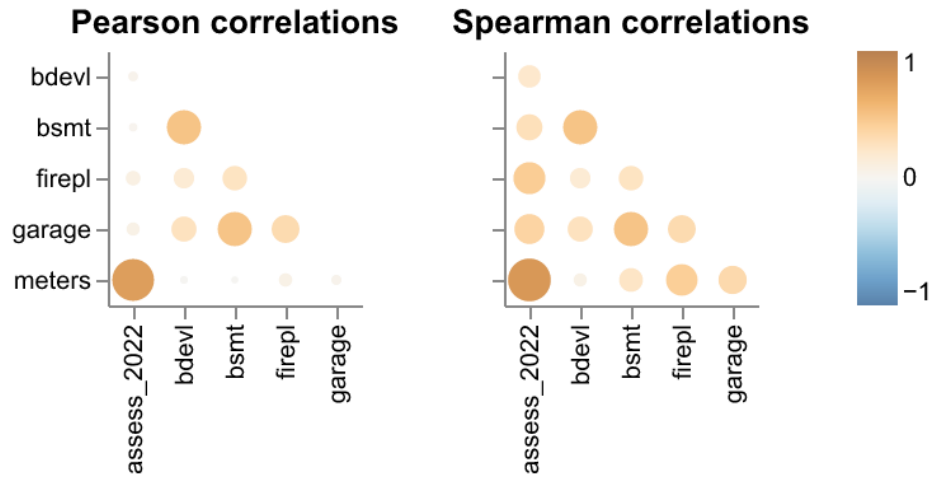


Figure 3: Correlation plot showing relationships between all features

As depicted in Figure 4, the combined visualizations provide valuable insights into the distribution and relationship between property sizes (meters) and house assessment values (`assess_2022`). The scatter plot showcases how larger properties tend to have higher assessment values, while the histogram and KDE plot reveal the overall distribution and density of assessment values in the dataset.

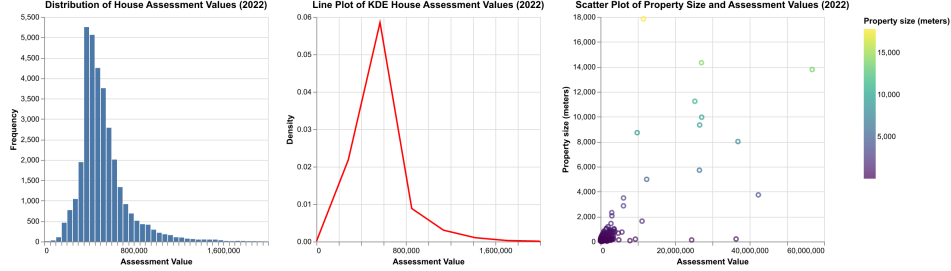


Figure 4: Distribution and relationships of house assessment values (2022), showing scatter plot, histogram, and kernel density estimation.

According to the scatter plot, the vast majority of our data points are concentrated within the $0 < \text{meters} < 2000$ range, with some more within the $2000 < \text{meters} < 5000$ range, and there are a few outliers in the $\text{meters} > 5000$ range that make it difficult to get a closer look on the majority of our data. We can see a similar story on the `assess_2022` feature as it has the vast majority of points in the low-mid range and some outlying extreme values.

We have split the `housing_df` DataFrame into training and test datasets using an 70-30 `split`, categorized features into categorical features (e.g., `garage`, `firepl`, `bsmt`, `bdevl`) and numeric features (e.g., `meters`), applied one-hot encoding transformation for categorical features and `StandardScaler` transformation for numeric features by using `Scikit-learn` (2024). We then created a column transformer `preprocessor` through combining these transformations to apply to the dataset and visualizes the `preprocessor`.

We performed analysis to identify any potential issues with the relationships between features and the target, as well as between the features themselves. By detecting anomalous correlations, we can uncover redundant or highly correlated features that might lead to multicollinearity, which can negatively affect the performance of machine learning models. In particular, removing or adjusting features that are too strongly correlated ensures that the model can learn meaningful, independent relationships, improving its generalization ability and interpretability. Additionally, examining correlations helps us assess if certain features are disproportionately influencing the target variable, which could indicate bias or overfitting.

During the validation step, we ensured the correctness of categorical feature levels by checking for string mismatches or unexpected single values. All categorical columns in the training dataset (`garage`, `firepl`, `bsmt`, and `bdevl`) passed validation with the expected levels of 'Y' and 'N'. Additionally, we analyzed correlations between the target variable and features, as well as correlations between features themselves. No anomalous correlations were detected in

either case. For correlation analysis, we excluded pairs of binary features that are inherently correlated due to their inverse relationship, such as `Garage_Y` and `Garage_N`, where one feature is always the inverse of the other. These relationships are expected and not flagged as anomalies. By excluding such pairs, we ensured the analysis focuses only on truly anomalous correlations, avoiding false positives while identifying potential issues in the dataset.

5. Train, cross validate and evaluate a Ridge regression model

We have split the features and target variable (`assess_2022`) into training and testing datasets. Then, we have created a pipeline combining the column transformer (`preprocessor`) and the Ridge Regression model (Hoerl and Kennard (1970)). Using 5-fold cross-validation on the training data, we evaluated the pipeline on multiple metrics and computes train and validation scores. Finally, it outputs the aggregated train and validation scores to assess the model's performance (Learn (2024)). Table 3 summarizes the cross-validation results, including the mean and standard deviation of train and validation scores across 5 folds. For our baseline, we created simple dummy model that uses mean values for predictions. Table 2 summarizes the baseline's cross-validation results

Table 2: Baseline cross-validation results showing the mean and standard deviation for train and validation scores across 5 folds.

	Unnamed: 0	mean	std
0	fit_time	0.013	0.002
1	score_time	0.004	0.001
2	test_score	-0.003	0.005
3	train_score	0	0

By comparing our model's performance in Table 3 with our baseline's performance in Table 2, we can see a significant improvement in scores which supports our model's validity.

Table 3: Cross-validation results showing the mean and standard deviation for train and validation scores across 5 folds.

	Unnamed: 0	mean	std
0	fit_time	0.024	0.011
1	score_time	0.004	0.001
2	test_score	0.564	0.235
3	train_score	0.575	0.078

We fit the pipeline on the training dataset (`X_train` and `y_train`) to train the Ridge Regression model (Hoerl and Kennard (1970)). Then evaluates the trained pipeline on the test dataset

(`X_test` and `y_test`) which calculates the R^2 (coefficient of determination) to measure how well the model explains the variance in the test data.

6. Predict housing prices with new data

We created a Pandas DataFrame containing information about 10 houses that we wish to predict the value of. Table 4 summarizes the attributes of ten houses, including property size and various categorical features, used for predictions.

Table 4: Attributes of ten houses used for prediction.

	meters	garage	firepl	bsmt	bdevl
0	174.23	Y	Y	Y	N
1	132.76	Y	N	Y	Y
2	90.82	Y	N	N	Y
3	68.54	N	N	N	N
4	221.3	Y	Y	Y	Y
5	145.03	N	N	N	Y
6	102.96	N	N	Y	Y
7	164.28	Y	Y	N	N
8	142.79	N	Y	Y	N
9	115.94	Y	N	Y	Y

We applied the trained pipeline to predict housing prices values based on the features in the DataFrame containing new data Hastie, Tibshirani, and Friedman (2009). The predictions are stored in a new pandas DataFrame.

We combined the original features from the new data and the predicted values into a new pandas DataFrame. The new pandas DataFrame provides an overview of the predictions. Table 5 displays the predicted house values for the ten houses based on their attributes, including property size and categorical features.

Table 5: Predicted house values for the ten houses based on their attributes.

	meters	garage	firepl	bsmt	bdevl	Predicted_Values
0	174.23	Y	Y	Y	N	537035
1	132.76	Y	N	Y	Y	457895
2	90.82	Y	N	N	Y	360993
3	68.54	N	N	N	N	212530
4	221.3	Y	Y	Y	Y	700375
5	145.03	N	N	N	Y	445199

Table 5: Predicted house values for the ten houses based on their attributes.

	meters	garage	firepl	bsmt	bdevl	Predicted_Values
6	102.96	N	N	Y	Y	344127
7	164.28	Y	Y	N	N	515519
8	142.79	N	Y	Y	N	419402
9	115.94	Y	N	Y	Y	418258

7. Visualization for predictions

We have created line charts to visualize the relationship between property size (**meters**) and predicted housing prices (**Predicted_Values**), colored by different categorical features. This visualization highlights how categorical features interact with property size to influence predicted house values“Altair Tutorial, Exploratory Data Visualization with Altair” (2024). As shown in Figure 5, property size has a significant impact on house value predictions, with notable variations across categorical features.

Discussion

The Ridge regression model demonstrates reasonable performance, with a training score of 0.575 and a test score of 0.564, suggesting it generalizes well without significant overfitting. However, the large variability in the test score (± 0.235) points to potential instability, which could be improved by tuning the model’s hyperparameters. The model is computationally efficient, with very low fit time (0.01s) and score time (0.003s), making it suitable for rapid training and prediction. While the overall performance is acceptable, further improvements could be made through hyperparameter optimization or exploring alternative models.

Our findings suggest that the aforementioned 10 houses are expected to be valued at: Table 6

Table 6: Predicted house values

	Predicted_Values
0	537035
1	457895
2	360993
3	212530
4	700375
5	445199
6	344127
7	515519
8	419402

Table 6: Predicted house values

	Predicted_Values
9	418258

We have also noticed that there is a correlation between a house’s price and its property’s size, whether or not it has a garage, fireplace, basement, and whether or not the building has been evaluated. This is consistent with our expectations as the larger a property is, and the more features it has (basement, garage, etc.) the higher its value becomesFastExpert (2024). This also raises further questions such as what other features of a house or property affect its value? characteristics for future consideration include: Number of bedrooms, number of bathrooms, indoor space, outdoor space, and number of floors.

Limitations

While the model performs well, several limitations need to be addressed:

1. Assumptions of Ridge Regression: Ridge regression assumes linear relationships between features and the target variable. While effective in addressing multicollinearity, it may fail to capture complex, non-linear interactions in the data. The regularization term could over-penalize certain features, potentially underestimating their influence.
2. Dataset Assumptions: The analysis assumes independence between observations, which might not hold if regional trends or external factors influence house prices. While pre-processing addressed issues like missing values, duplicates, and outliers, any inaccuracies or biases in the original data could impact the model’s predictions.
3. Feature Selection: The model is limited to the selected features (meters, garage, firepl, bsmt, bdevl), which do not comprehensively represent all factors affecting house prices. Features like the number of bedrooms, bathrooms, or outdoor space could improve the model’s performance.

Future Work

To address these limitations, future iterations could incorporate non-linear models such as decision trees or ensemble methods to capture more complex relationships. Including additional features like the number of bedrooms, bathrooms, or proximity to amenities could further refine the predictions. Spatial modeling or time-series analysis might also better account for dependencies in the data.

By acknowledging these limitations and proposing improvements, this project lays the groundwork for future enhancements in house price prediction.

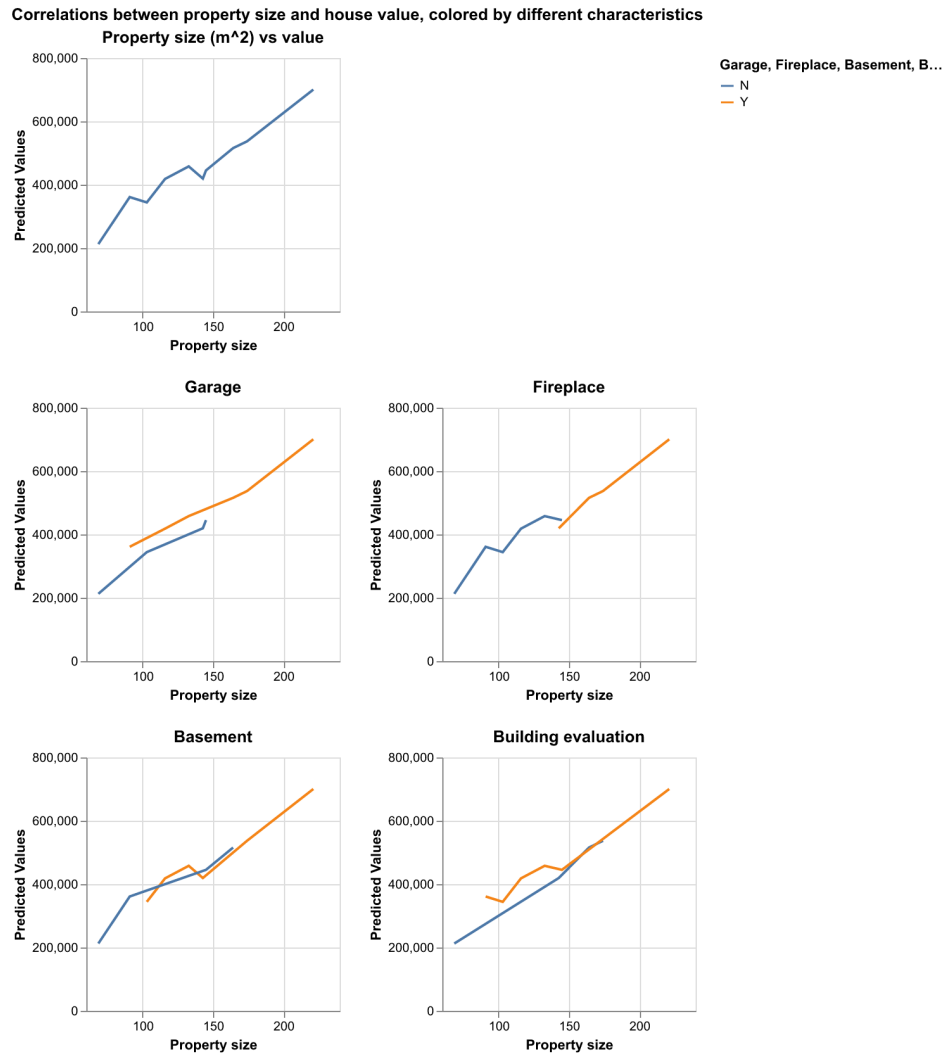


Figure 5: Line charts showing correlations between property size and house value predictions, colored by various categorical features.

References

- “Altair Tutorial, Exploratory Data Visualization with Altair.” 2024. 2024. <https://altair-viz.github.io/altair-tutorial/README.html>.
- County, Strathcona. 2023. “2023 Property Tax Assessment.” <https://www.strathcona.ca/services/assessment/>.
- FastExpert. 2024. “How Much Value Does a Garage Add to a House?” 2024. <https://www.fastexpert.com/blog/how-much-value-does-a-garage-add-to-a-house/>.
- Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- Hoerl, Arthur E., and Robert W. Kennard. 1970. “Ridge Regression: Biased Estimation for Nonorthogonal Problems.” *Technometrics* 12 (1): 55–67.
- Learn, Scikit. 2024. “Model Selection and Evaluation.” 2024. https://scikit-learn.org/stable/model_selection.html.
- Team, Quarto Development. 2024. “Quarto: An Open-Source Scientific and Technical Publishing System.” <https://quarto.org/>.