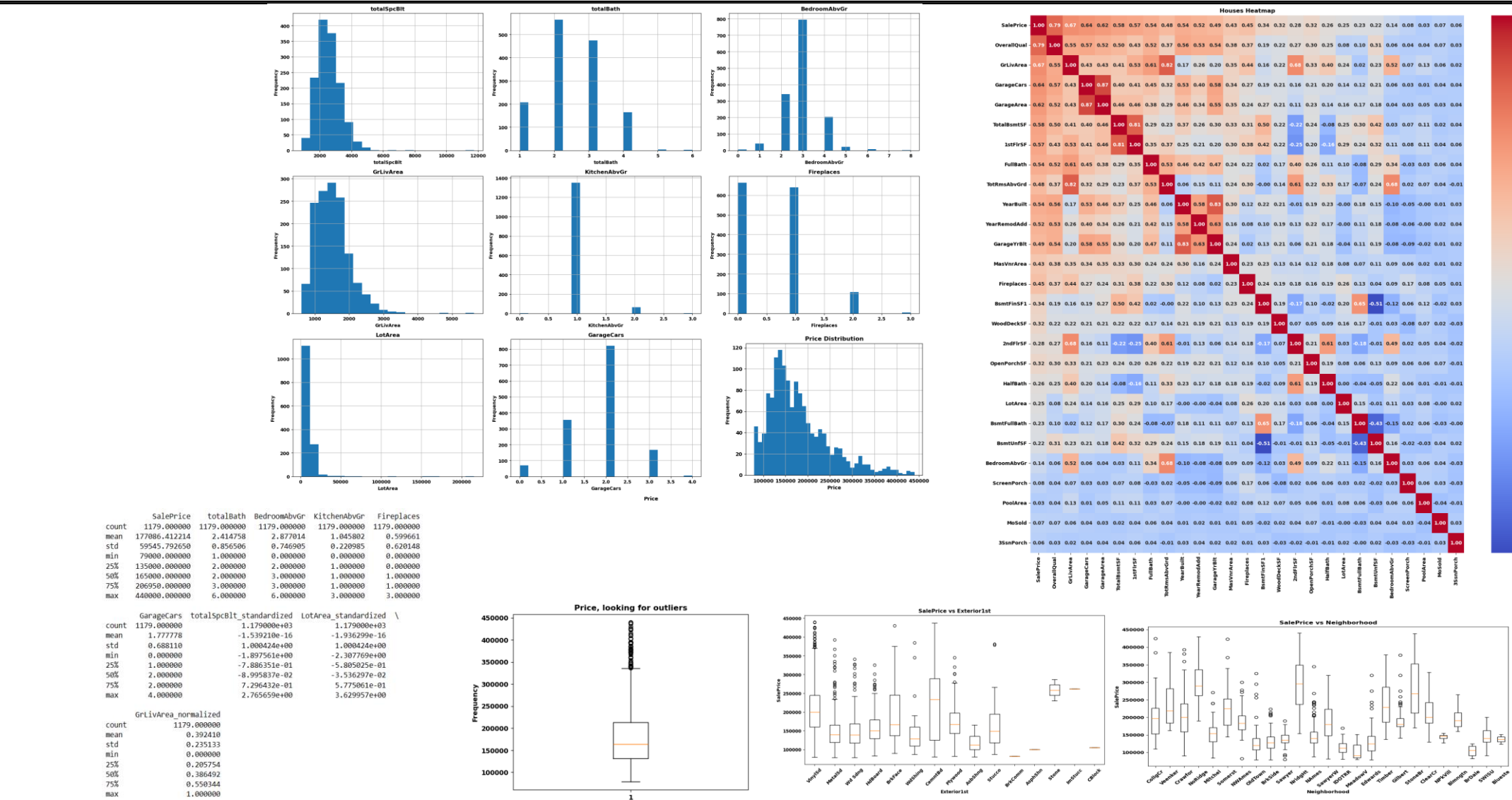


Description and Motivation

- Building and comparing two machine learning models: Lineal Regression vs Decision Tree Regression.
- This piece of work aims to clarify if houses prices can be accurately predicted by the most basic variables that every house should have such as kitchen, living room space, total space built, etc.
- The process followed to build the models was kept as simple as possible. When data is limited, simpler models are many times more effective ¹.
- Common sense suggests that houses prices should be predictable based on the kind of the aforementioned features.
- Since I wanted to find my very own discoveries, only few references were briefly read before starting, and during the project. Most of the text relay on my own experience/analysis. There isn't many specific references, although used them as general guide in the whole process .

Initial Analysis of Dataset

- The chosen data set for the project is the “Ames, Iowa House Market” compiled by Dean De Cock ².
- The original data set is form by features that every house could have (pool, kitchen, living room, lot area, etc).
- The original data set has1460 row x 81 columns (Including “Id”). Some columns were deleted since they had many empty values. After data cleaning, 1179 records x 36 columns were left to do the analysis.
- There were many outliers across the different variables and as example we can take a look at the target variable: “SalePrice”.
- Some of categorical variables seem to have a high variability. For instance, sale price-neighbourhood shows a significant variation in the mean of the house prices this means that the location definitely affect to the price. However, others also shown similar behaviour but after further examination the percentage of the houses affected by those variables in their price were quite small is the case of the sale price vs exterior1st (exterior covering the house).
- As was diving into the data set, the former motivation changed from “trying to build two models to make accurate predictions, then compare them” to focus on the basic features that every basic house should have; lot area, bathroom, garage, kitchen, and so on as well as others relevant such as neighbourhood. I understood that to build better models I should start by the basics.
- Applied feature engineering to some of the features(Standardization, normalization, one-hot encoding). Created new variables combining others. The feature totalSpcBlt_standardized is a standardization of the figures for the total space built (surface built in first floor + surface built second floor, etc).
- Looking at the correlation heatmap we can easily spot the highest correlated variables against the dependent variable “SalePrice”. It is worthy to mention the variables GrLivArea (ground living area) coefficient of 0.67 or “GarageCars” with a coefficient of 0.64. Even though OverallQual (overall material and finish quality) has the strongest correlation with a coefficient of 0.79 has not been included since the focus of the study is “tangible house parts”.
- Examined the data distribution to determine its compatibility with the models. Data distribution of the Sale Price, total space built, and ground living area is right-skewed whereas lot area follow a normal distribution.
- Multicollinearity was checked by applying VIF (Variance Inflation Factor), many columns had multicollinearity. High multicollinearity may make the model unstable therefore, some of the columns were deleted ³.



Linear Regression

- Linear Regression is a model that assumes a relationship between the dependent variable and the independent variable.
- Linear Regression predict the dependent variable based on a function of the independent variables. It uses a linear equation $Y = a + b \cdot X$ where “a” is the intercept and “b” is the slope of the lineal representation.
- The coefficients measure the effect of each independent variable on the dependent var iable. They are denoted by the Greek letter β . The intercept is represented by β_0 and indicates the value of “Y” when $x = 0$. Whereas β_1 represents the slope, the slopes indicates how much changes “Y” for each unit of β_1 .

Pros

- Simplicity: Linear regression is easy to grasp, therefore, is very appropriate for situations where transparency is crucial.
- Computational efficiency. The model require from little physical resources this means is an agile model even when training large datasets.
- Suitability for continuous data, this is, Lineal Regression performs well with continuous data and for our case (house prices) is an idyllic option.
- It has the potential to be highly accurate when the relationship between variables is linear.

Cons

- Sensitivity to outliers. In this sense, W. Choi highlight the impact of outliers on linear regression models.
- Assumes linearity. This is definitely not the real-world scenario where sometimes there isn't any relationship between variables.
- Multicollinearity. The model takes for granted that there is not high correlation between the independent variables, these presumption are wrong and can lead to make unstable estimations of the model.
- There is a high risk of overfitting specially in cases where the model includes multiple independent variables.

Decision Tree Regression

- Decision Tree is a non-parametric supervised algorithm, it means that they don't need a predetermined set of parameters to start the training.
- The model is similar to an upside-down tree. It consists of a root node that represent the entire population, decision nodes that split the data into smaller subsets until the terminal node. Leaf (or terminal) nodes represent a class label in the data according to the path taken.
- It tries to split the data into smaller groups based on the features, at the time also try to make as similar as possible in terms of the target value. The final prediction is based on the terminal nodes.

Pros

- Easy to understand and interpret, the decisions can be seen through the diagram tree.
- Can handle both numerical and categorical data. Can address multi-output problems.
- They do automatic feature selection by identifying the best features for splitting data
- Can manage big amounts of data sets without the need to adjust the range of different features before using them.

Cons

- Decision Trees are prone to overfitting, especially when the tree is deep. This fact can make the model sensitive to fluctuations.
- Small variations in the data can lead to the generations of different trees making them not that reliable for making predictions.
- Simple decision trees may not have the same capacity to make predictions as models are more complex this can lead to results with the presence of bias.
- Decision Trees have high variance and can create too many complex trees that don't fit well in new data.

Hypothesis Statement

- Linear Regression model should perform better since it is supposed that the selected features have a strong relationship with the target variable, however in contrast with Decision Tree Regression, linear regression might not capture complex non-linear relationships therefore the performance of both models will probably be similar.
- Overall, we expect that the Linear Regression model will outperform the Decision Tree Regression model ⁴
- Perhaps seems too obvious but errors in the prediction of house prices will likely be more significant in houses that deviate from the average price.

Methodology

- Set a seed for reproducibility of the two models.
- “Holdout validation technique” was used. Data was split into 70% training, 30% test.
- Built a simple Linear Regression model considering all the original variables. Examining model performance gave me a general insight into the whole (models, features, etc).
- Based on my years of experience in the construction field, I proceeded to make a first selection of features, built another model with these features and examined the results.
- Feature selection. Applied Lasso regression for variable selection and regularization, it helped to identify and select the most relevant features for the final model.
- Built another linear regression model, examined the results.
- Continued with feature selection, applied VIF(Variable Inflation Factor) to check multicollinearity. Based on the results, deleted highly correlated features to avoid overfitting or bias in the models.
- Built different versions of the two models and checked results (R-squared, Root Mean Squared Error, Mean Absolute Error, Residuals, charts results for the two models). Selected final parameters such as “ maximum number of splits”.

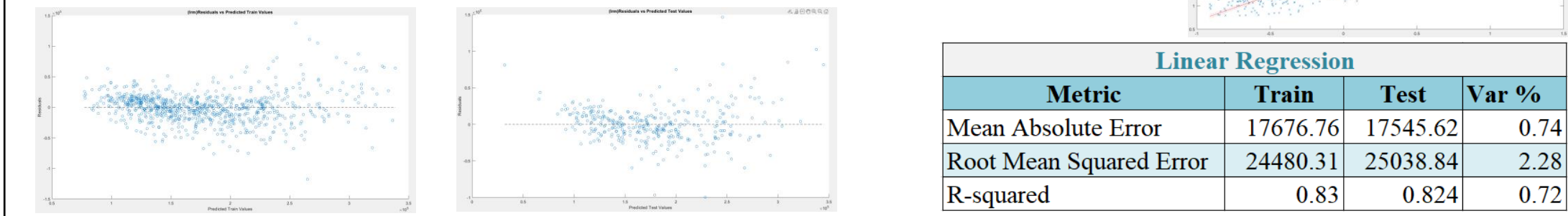
Choice Of Parameters And Experimental Results

Common points for both models

- Feature selection for both models was initially performed discretionarily, followed by the application of Lasso regression and the Variance Inflation Factor(VIF) analysis.
- Initial results before feature selection were unstable and performance was quite poor. The residual errors were over 100,000 units for example.
- Plot the residuals Vs predicted train values, and residuals vs predicted test values for further understanding of the metric. The charts for each model are below.
- Calculated the formulas for MAE, RMSE through coding in MATLAB.
- The biggest difference between models happens in the test. Right side upper corner.

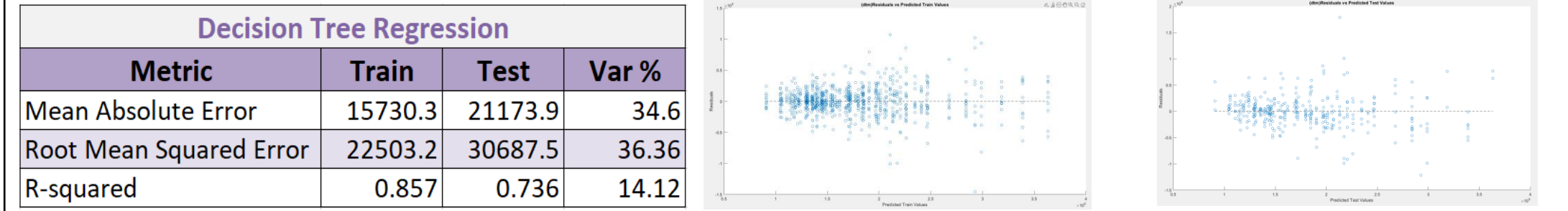
Linear Regression

- On the right side, the chart for the final Lineal Regression model before any evaluation against training/test models. this is, predictors Vs target value.
- Parameters were automatically defined by MATLAB. “fitlm” function estimates the parameters for each predictor by minimizing the sum of the squares of the residuals.



Decision Tree Regression

- On the right side, the decision tree regression chart, before any evaluation against training/test, shows the root node based on the total space built “totalSpcBlt”. The final values at the leaves (on the right side) represent the predictions
- Parameters were chosen manually using the Fibonacci series. After experimenting with different options, selected the combination of parameters with the best outcomes. For some combinations of parameters, the outcomes were considerably bad.



REFERENCES

[1] P. Mehta et al., ‘A High-bias, low-variance introduction to Machine Learning for physicists’, Physics Reports, vol. 810, pp. 1–124, May 2019, doi: 10.1016/j.physrep.2019.03.001.

[2] D. De Cock, ‘Ames, Iowa: Alternative to the Boston Housing Data as an End of Semester Regression Project’, Journal of Statistics Education, vol. 19, no. 3, p. 8, Nov. 2011, doi: 10.1080/10691898.2011.11889627.

[3] J. I. Daoud, ‘Multicollinearity and Regression Analysis’, J. Phys.: Conf. Ser., vol. 949, p. 012009, Dec. 2017, doi: 10.1088/1742-6596/949/1/012009.

[4] P. Jadhav, V. Patil, and S. Gore, ‘A Comparative Study of Linear Regression and Regression Tree’, SSRN Journal, 2020, doi: 10.2139/ssrn.3645883.

[5] Z. Zhang, ‘Decision Trees for Objective House Price Prediction’, in 2021 3rd International Conference on Machine Learning, Big Data and Business Intelligence (MLDBDI),

Taiyuan, China: IEEE, Dec. 2021, pp. 280–283, doi: 10.1109/MLDBDI54094.2021.00059.

[6] J. R. Quinlan, ‘Induction of decision trees’, Mach Learn, vol. 1, no. 1, pp. 81–106, Mar. 1986, doi: 10.1007/BF00116251.

[7] R. P. Ribeiro and N. Moniz, ‘Imbalanced regression and extreme value prediction’, Mach Learn, vol. 109, no. 9–10, pp. 1803–1835, Sep. 2020, doi: 10.1007/s10994-020-05900-9.

[8] H. Yu and J. Wu, ‘Real Estate Price Prediction with Regression and Classification CS 229 Autumn 2016 Project Final Report.’ Available: https://cs229.stanford.edu/proj2016/report/WuYu_HousingPrice_report.pdf

[9] S.-W. Choi, ‘The Effect of Outliers on Regression Analysis: Regime Type and Foreign Direct Investment’, QJPS, vol. 4, no. 2, pp. 153–165, Jun. 2009, doi: 10.1561/100.00008021.

[10] [10] “MATLAB Documentation - MathWorks United Kingdom.” https://uk.mathworks.com/help/matlab/index.html?tid=hc_panel