DATA MINING AND MACHINE LEARNING

ASSESSMENT 3

ANTHONY OGUNNA

2413689

An Assignment Submitted in the partial fulfilment for

the award of Master of Science degree in Data Analytics and Technologies

**Tutor :** Dr. Pradeep Hewage

April 2025

Table of Contents

Abstract ………………………………………………………………………………………iii

Keywords: Data mining, Machine Learning models, Modeling techniques, Data science project management, Models Evaluation.

Word count : 2133

ABSTRACT

This study focuses on using the regression modeling technique to predict the prices of house in a given location with machine learning predictive algorithms. For effective assessment of the project, the Cross Industry Standard Process for Data Mining (CRISP-DM) methodology is employed to build a comprehensive framework to manage the process from the business understanding phase down to the deployment stage.

# CHAPTER 1

## INTRODUCTION

### 1.1 Business understanding

Housing, which is one of the most basic needs for human existence, is at the center of economic firmness, along with food, water, and many other necessities (Zulkifley et al., 2020). House pricing, which deals with the valuation of residential properties, is influenced by factors such as supply and demand, population growth, employment rates, availability of mortgage financing and affordable interest rates, property size and features, economic conditions, government policies, and others. It plays an important role in shaping economic stability and overall societal wellbeing. In addition, an insight into house pricing is essential for buyers, sellers, investors, and policymakers, as it helps to make informed decisions regarding property transactions and market trends (Kang et al., 2021). Furthermore, the value of statistical models and machine learning in the profession of property valuation is a topic of continuous debate among appraisers and real estate market experts (Forys, 2022). The key objective of this project is to analyse and predict house pricing as this is one of the key business problems in real estate business using data-driven approaches with machine learning by developing predictive models and evaluate them to get the best model that meet the business need to help stakeholders make informed decisions with respect to house prices.

### 1.2 Data understanding

The housing price dataset contains the following columns:

PARCELNO: unique identifier for each property. About 1% appear multiple times.

SALE_PRC: sale price in U.S. dollars ($)

LND_SQFOOT: This is the total land area of the property measured in square feet

TOT_LVG_AREA: This is the total living (floor) area of the property measured in square feet

SPEC_FEAT_VAL: This is the estimated value of special features (e.g., swimming pools) in U.S. dollars ($)

RAIL_DIST: This is the distance from the property to the nearest rail line (an indicator of noise) measured in feet

OCEAN_DIST: This is the distance from the property to the ocean (feet)

WATER_DIST: This is the distance from the property to the nearest body of water eg lake, river etc measured in feet

CNTR_DIST: This is the distance from the property to the central business district measured in feet

SUBCNTR_DI: This is the distance from the property to the nearest subcenter (a secondary commercial hub) measured in feet

HWY_DIST:  This is the distance from the property to the nearest highway (an indicator of noise) measured in feet

age: This is the age of the structure in years

avno60plus: dummy variable for airplane noise exceeding an acceptable level with the threshold yes =1 and no = 0

structure_quality: This is a qualitative measure of the  quality of the property on numerical scale

month_sold: This is the month the property was sold in 2023 (1 = jan, 2 = feb, 3 = mar, etc)

LATITUDE: Geographical coordinate of the property

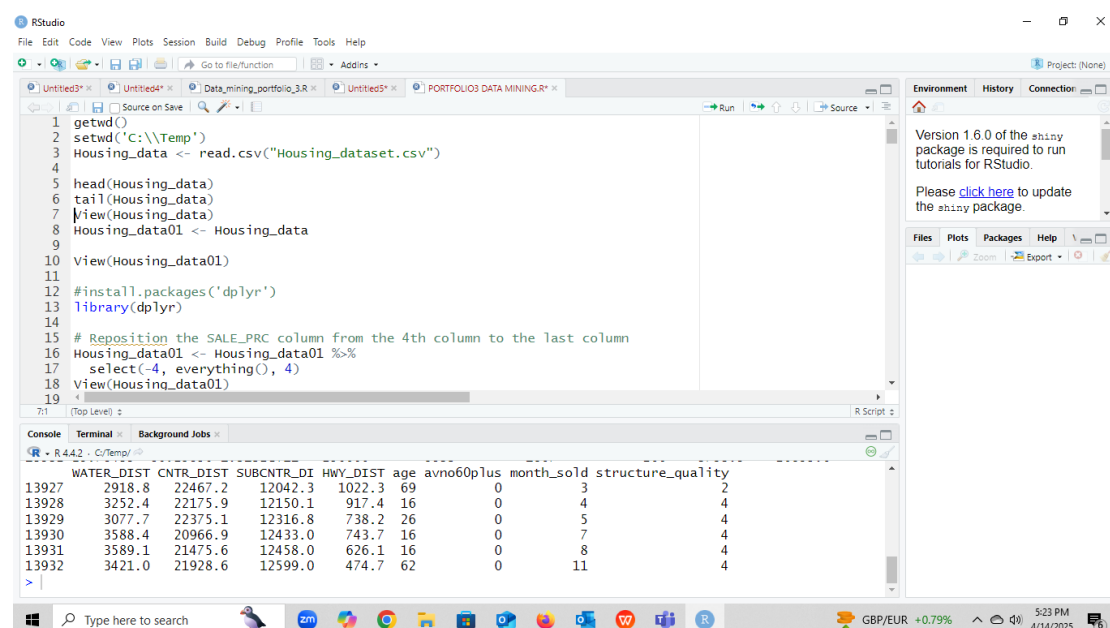LONGITUDE: Geographical coordinate of the property



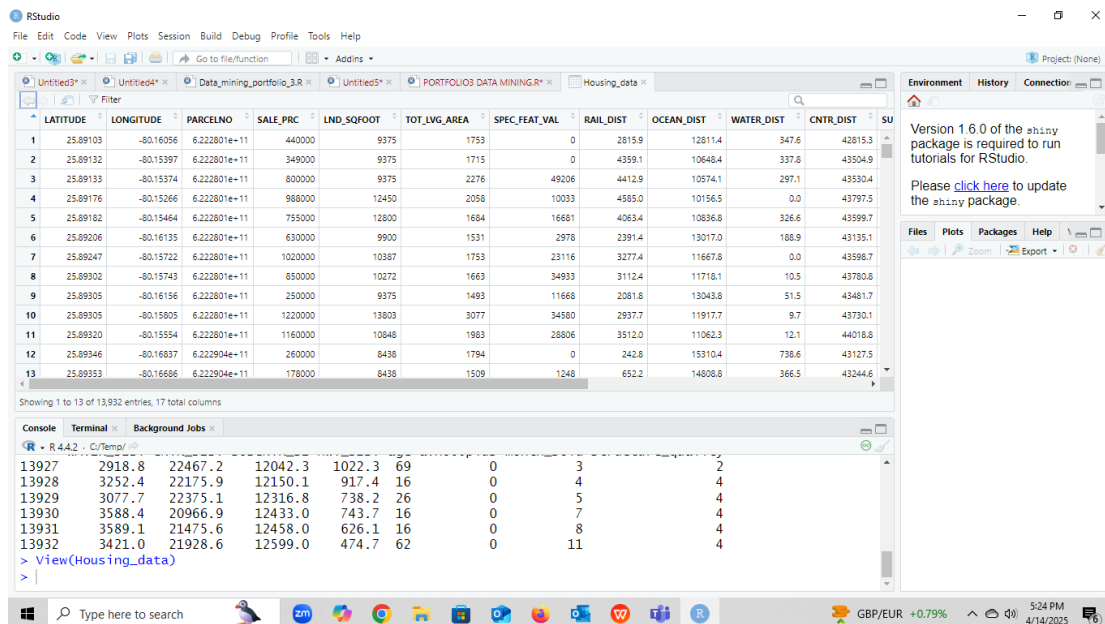Fig. 1: Loading the housing dataset on Rstudio

Fig. 2: A view of the housing dataset

## Check for variables with significant correlation with sale price

In the course of understanding the data, the researcher checked for the relationship of other variables with the house sales price in order to know which ones have significant impact on house sale price for model building using correlation analysis alongside heatmap visualisation in R as shown in fig.3, fig.4, fig.5 and fig 6 below. The outcome of the analysis shows that structure quality, HWY_DIST, SPEC_FEAT_VAL, TOT_LVG_AREA and LND_SQFOOT correlate positively with SALE_PRC while SUBCNTR_DI, CNTR_DIST, OCEAN_DIST, age, WATER_DIST and RAIL_DIST correlate negatively with SALE_PRC with LONGITUDE, LATITUDE, avno60plus, month sold and PARCELNO having insignificant impact on SALE_PRC since the houses are in the same location geographically and the other variables listed do not impact on house sale prices.
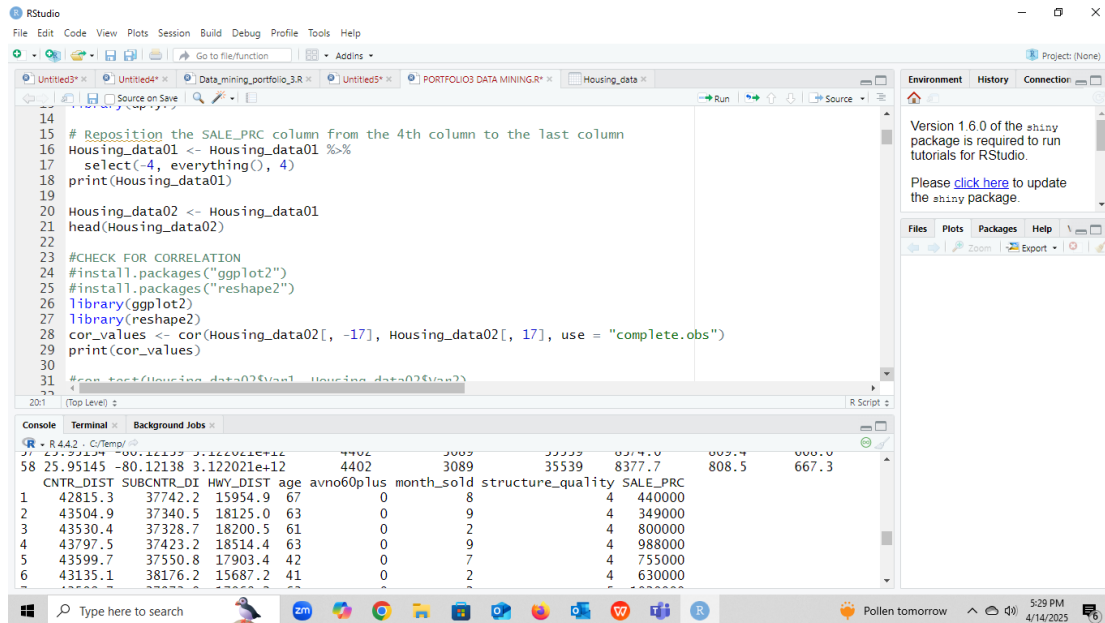
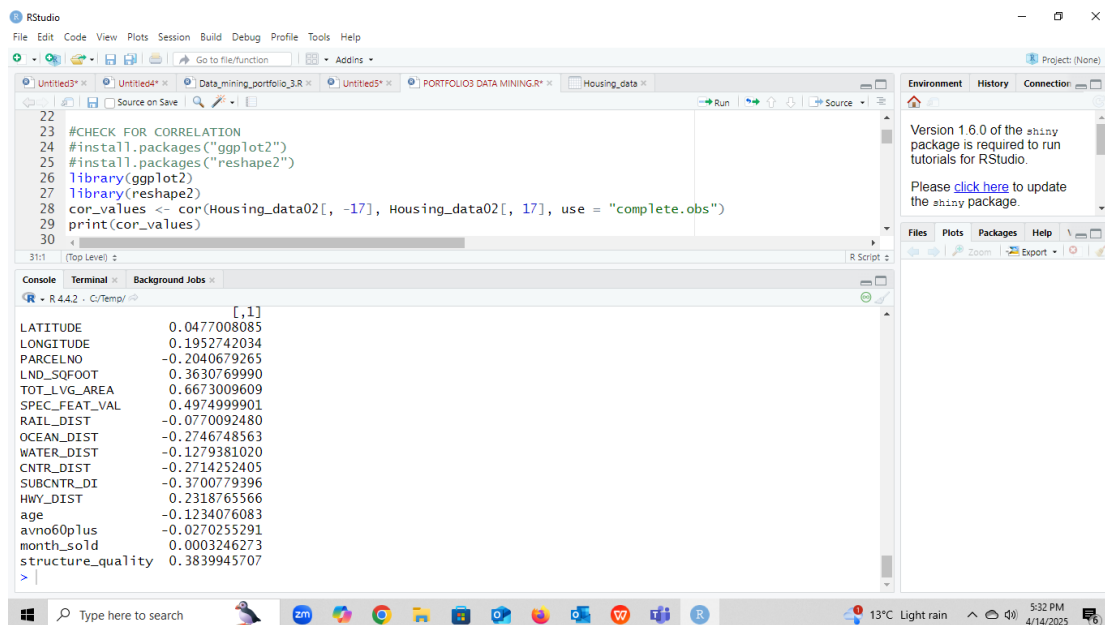Fig. 3: Reposition the SALE_PRC column from the 4<sup>th</sup> column to the last column



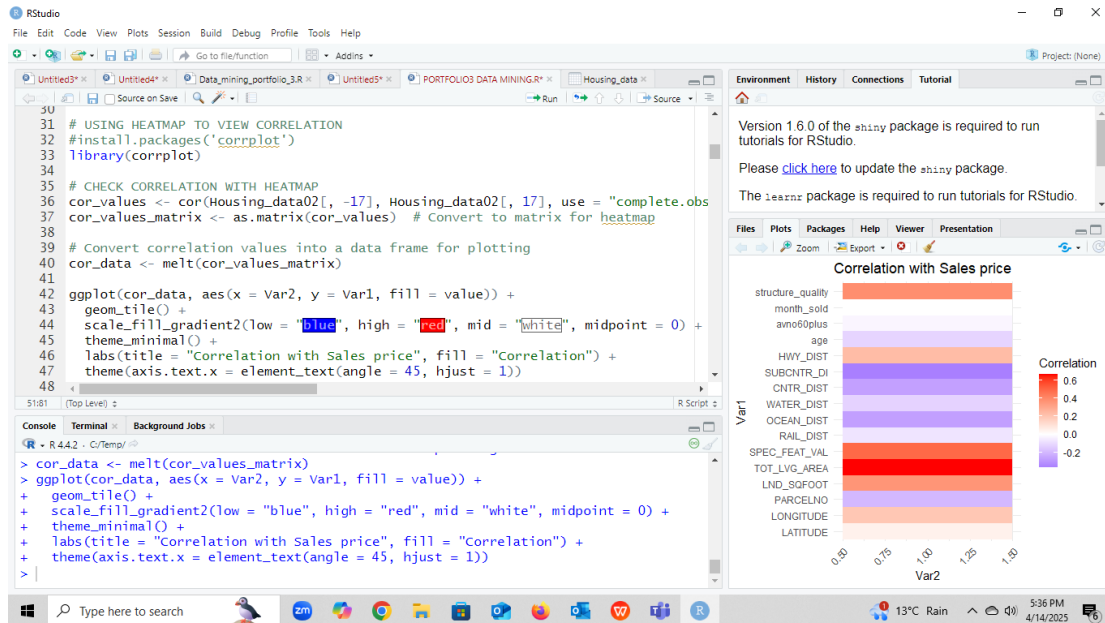Fig. 4: Check for the correlation of other variables with the SALE_PRC

Fig. 5: Correlation test on the relationship of SALE_PRC with other variables using Heatmap
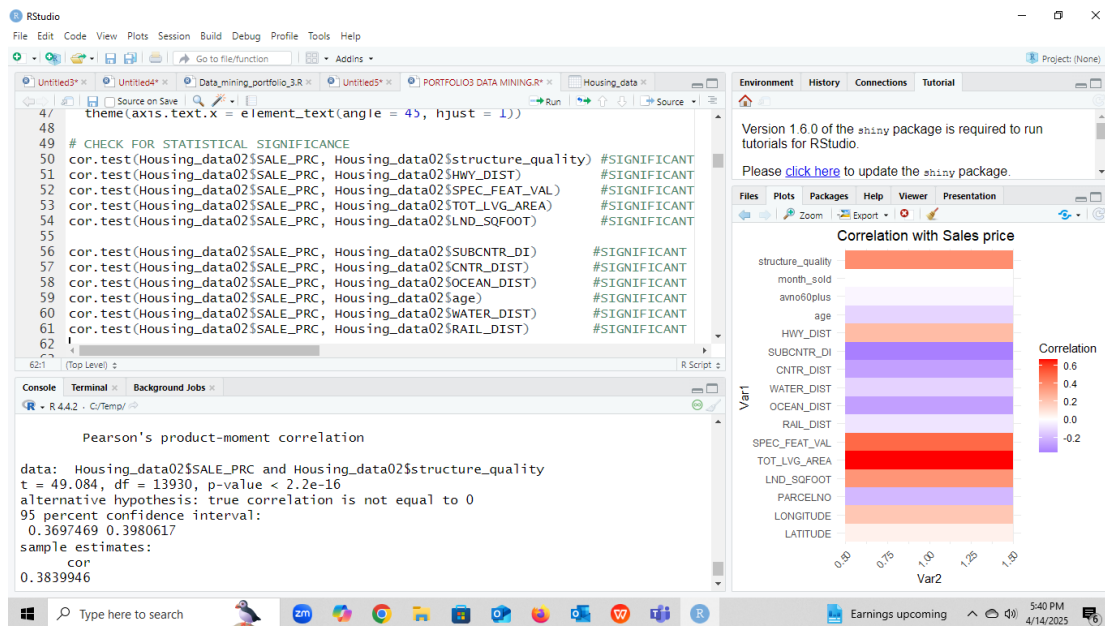


Fig. 6: Statistical test for correlation significance

# CHAPTER 2

## DATA PREPARATION

This phase of the task shows the in-depth data preparation process after identifying the variables that have significant impact on house sale price for model building. It ranges from check for duplicates , drop unwanted columns, check and handling of null values and outliers and carry out other necessary data cleaning procedures.

Below are screenshots of the data preparation process;

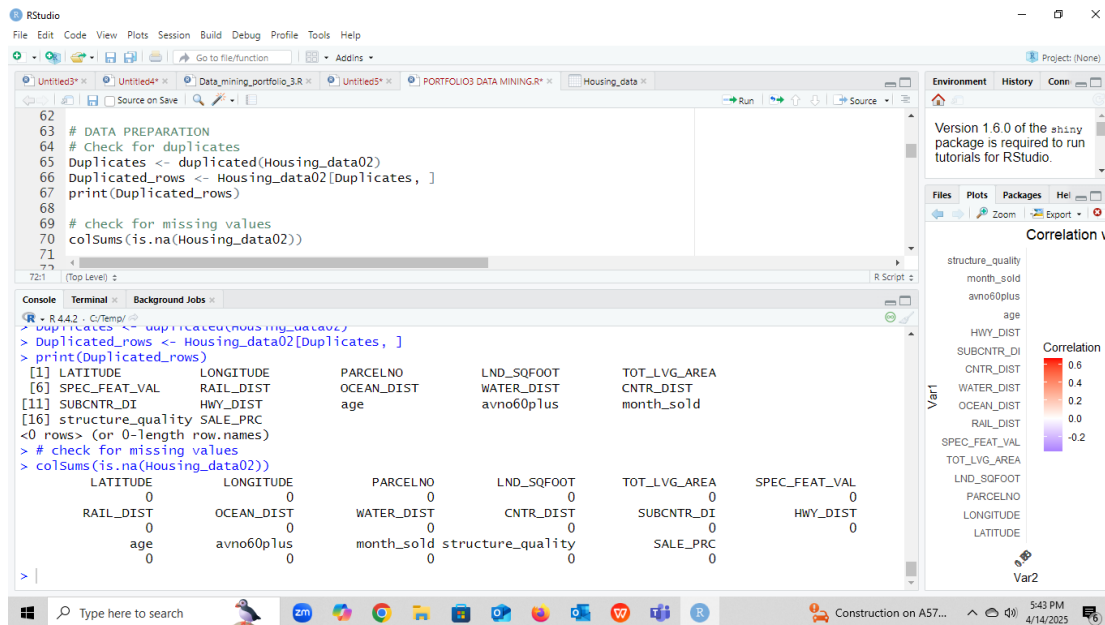## Check for duplicate



Fig. 7: Check for duplicates and null value

From the above screenshot in Fig. 7, it can be seen that the housing data does not contain any duplicated value after checking for duplicates using the duplicate function.

## Check for null value

The housing data does not contain null value as seen in Fig. 7 above.

## Drop columns with insignificant correlation on sale price

Fig. 8: Drop columns with insignificant correlation

After carrying out correlation analysis, the following columns; LATITUDE, LONGITUDE, PARCELNO, avno60plus and month_sold are considered to be less relevant in the model building for house price prediction due to their insignificant correlation with house sale price thereby removing them from the dataset as shown above in Fig.8

**Structure of Housing data**



Fig. 9: Structure of Housing data

The structure of the housing data for model building as shown in fig.9  has twelve (12) variables consisting of eleven explanatory variables and one dependent variable being

the sale price. The total number of observations is thirteen thousand nine hundred and thirty two with all columns data type being numerical.

## Handling outliers
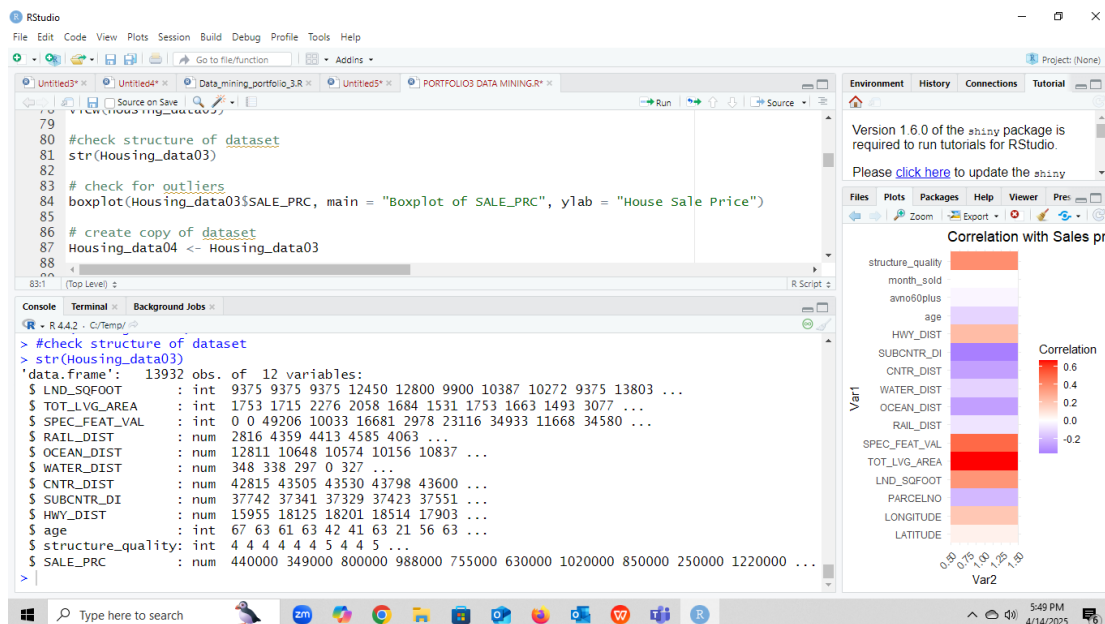


Fig. 10: Check for outliers

The above boxplot in Fig.10 shows the outliers chart. Because of the high margins in the house sale prices as a result of different components ranging from land and house sizes, special features among others, the house sale prices that appear as outliers but not truly outliers in the actual sense going by the context of the dataset will be included in the model building process for effective and accurate predictions.

CHAPTER 3

MODELING

## 3.1 Modeling technique

The algorithm to build and assess the various models for this project is the regression modeling technique. This technique model relationship between a dependent variable in this case the house sale price and other independent variables as contained in the house data as shown in Fig.8 as it gives solid predictions for continuous outcomes like the house sale price.

## 3.2 Machine Learning Models

In this project, eight (8) machine learning models are built, namely; multiple linear regression model, support vector regression linear model, support vector regression radial model, support vector regression polynomial model, decision tree model, random forest model_n100, random forest model_n200, and random forest model_n500 to predict the house sale price (SALE_PRC) using the following explanatory variables; LND_SQFOOT, TOT_LVG_AREA, SPEC_FEAT_VAL , RAIL_DIST, OCEAN_DIST, WATER_DIST, CNTR_DIST, SUBCNTR_DI, HWY_DIST, age and structure_quality. The model building process starts with setting seed for reproducibility at 123 to ensure that the same random results are obtained every time the code is executed as shown in Fig…… The data is then divided into training data and testing data with 70% of the total data accounting for the training data and the remaining 30% data assigned to the testing data. The next step in the model building phase is fitting the models on the training data where SALE_PRC is the dependent variable and the rest are the independent variables. When the model fitting phase is completed, prediction is then made on the testing data using the trained model as shown in fig. 12, fig. 13, fig. 14, fig. 15, fig. 16, fig. 17 and fig. 18 below . This process is carried out on the eight (8) models to train them with the training data before generating predictions using the testing data. At the end of this phase, the models are then evaluated based on their performance metrics to identify the model that best meet the business purpose for the project.
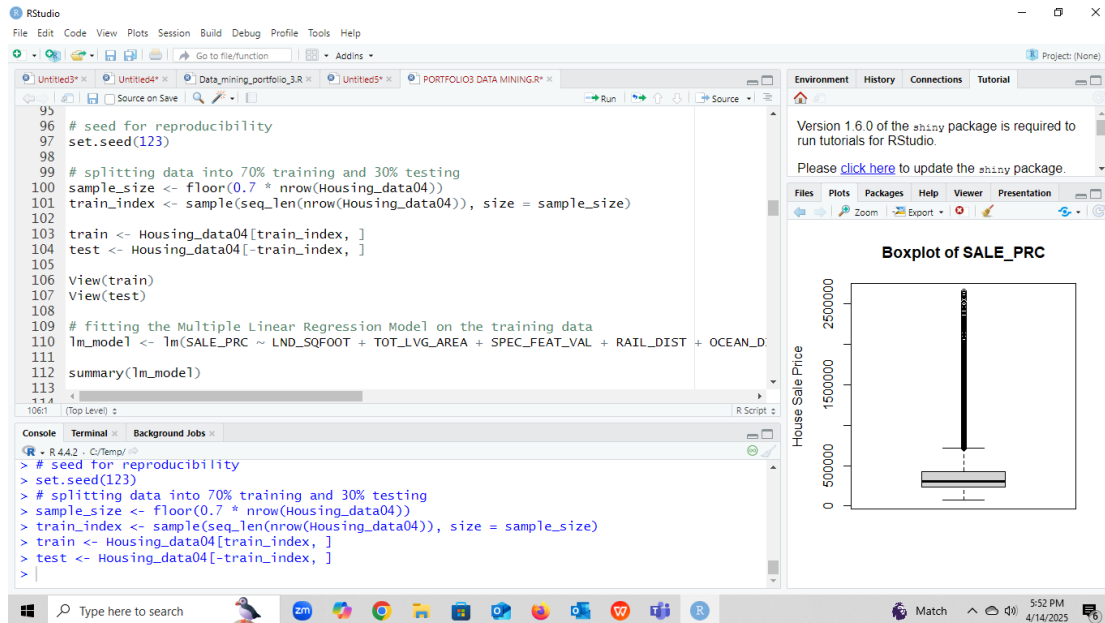
Fig. 11: Set seed and split data into training and testing data for model building



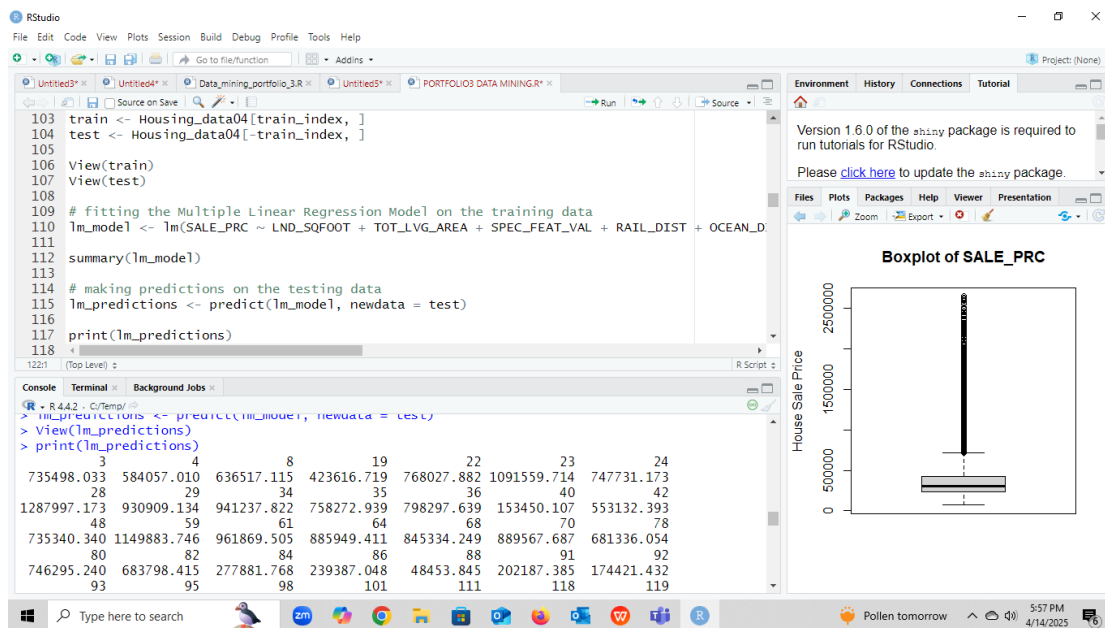Fig. 12: Fitting multiple linear regression model on training data and prediction

Fig. 13: Fitting support vector regression linear model on training data and prediction



Fig. 14: Fitting support vector regression radial model on training data and prediction

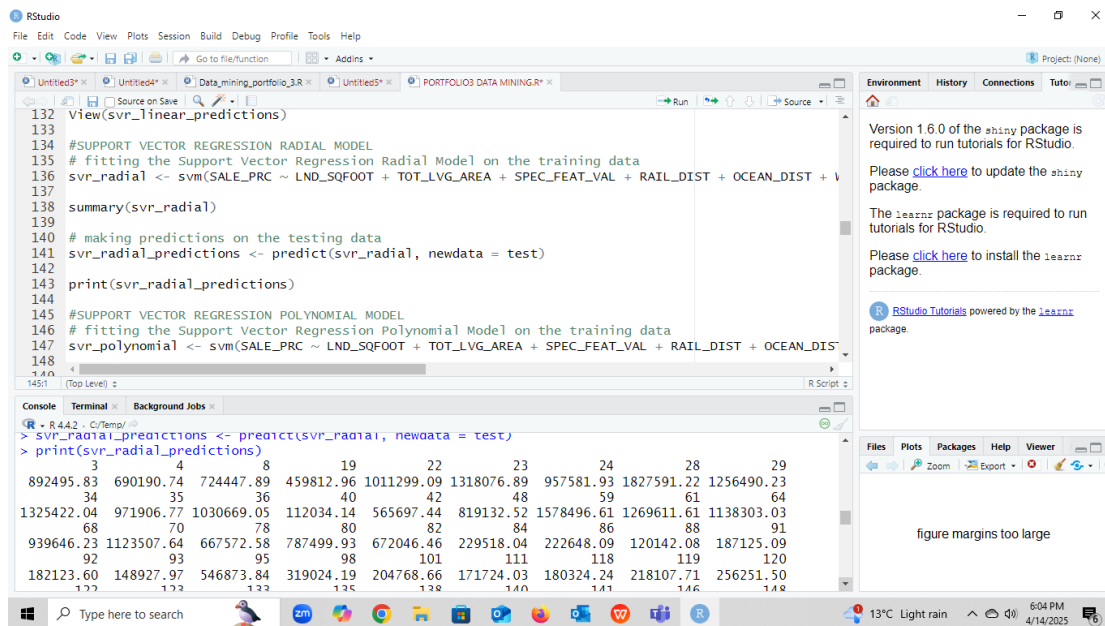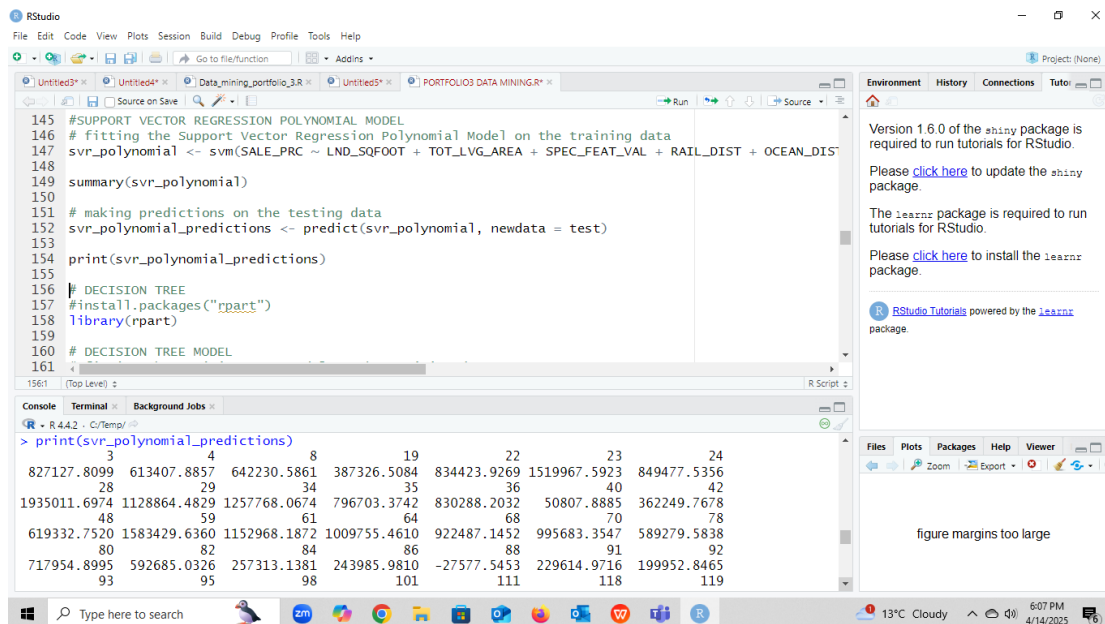Fig. 15: Fitting support vector regression polynomial model on training data and prediction
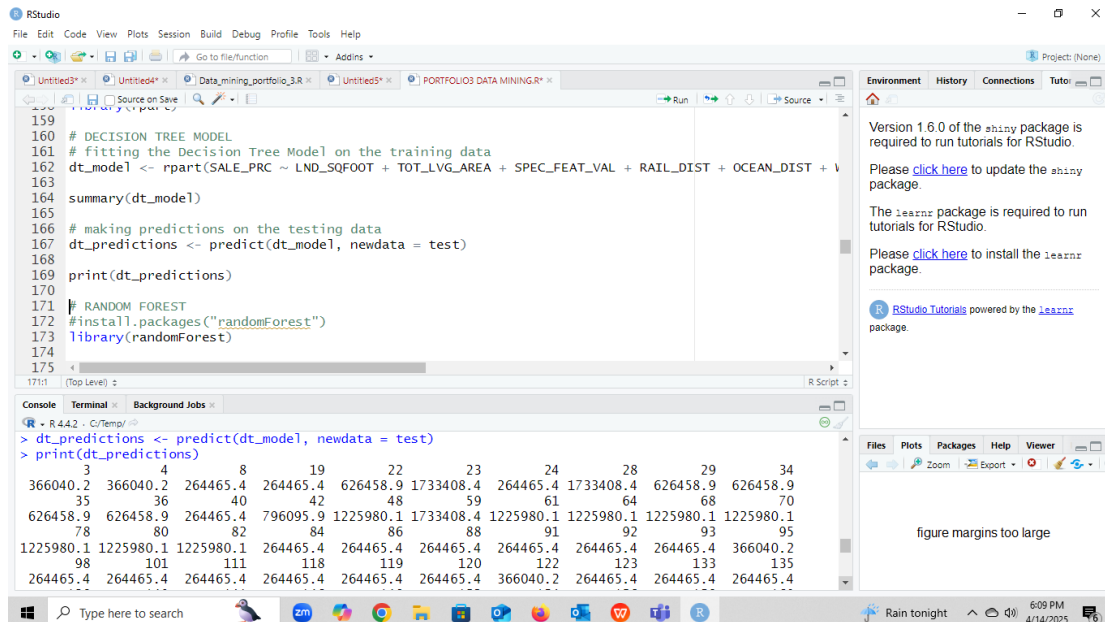


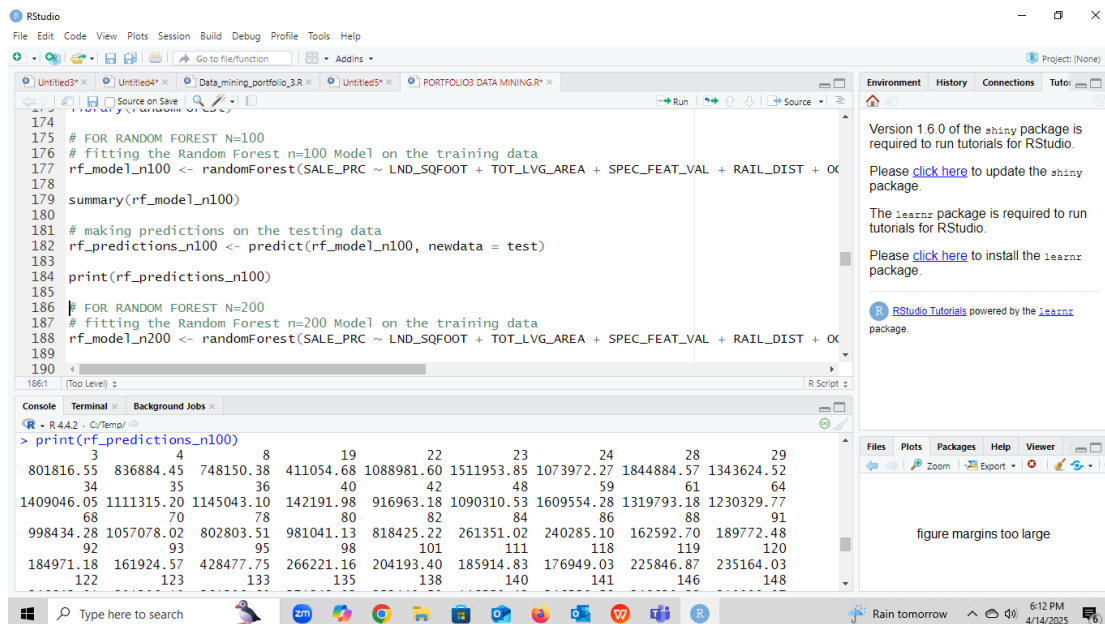Fig. 16: Fitting decision tree model on training data and prediction

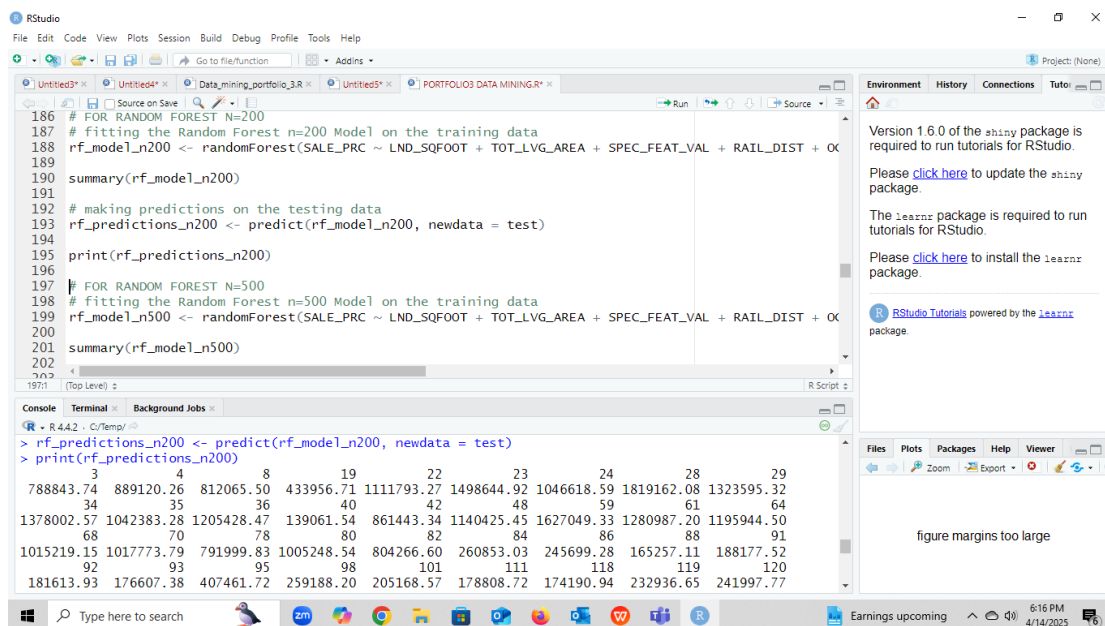Fig. 17: Fitting random forest model with n=100 on training data and prediction



Fig. 18: Fitting random forest model with n=200 and n=500 on training data and prediction

# CHAPTER 4

## EVALUATION  OF MODELS

The models are evaluated by calculating their individual performance metrics and compare using visualisation procedure to identify the best model with the least error using the mean squared error. This is done by using the postResample function from caret package to compute standard performance metrics for the regression models as shown in fig.19 and fig.20 below. The outcome from the evaluation phase shows the various models with their corresponding mean squared error values represented graphically as seen in the bar chart. From the mean squared error values obtained, support vector regression linear model has the least error making it the best fit for the business project. In order to optimise the support vector regression with linear kernel, hyper parameter tuning is carried out on the model to enhance its effectiveness and prediction capacity.

## Optimal model test

Given the following information   LND_SQFOOT: 11247, TOTLVGAREA: 4552, SPECFEATVAL: 2105, RAIL_DIST: 4871.9, OCEAN_DIST: 18507.2, WATER_DIST:375.8, CNTR_DIST: 43897.9, SUBCNTR_DI: 40115.7, HWY_DIST: 41917.1, age: 42 and structure_quality: 5 to make prediction on the house sale price using the optimised support vector regression model with linear kernel which is the optimal model in the project, the predicted house sale price is 837089.31 as shown in fig. 22.
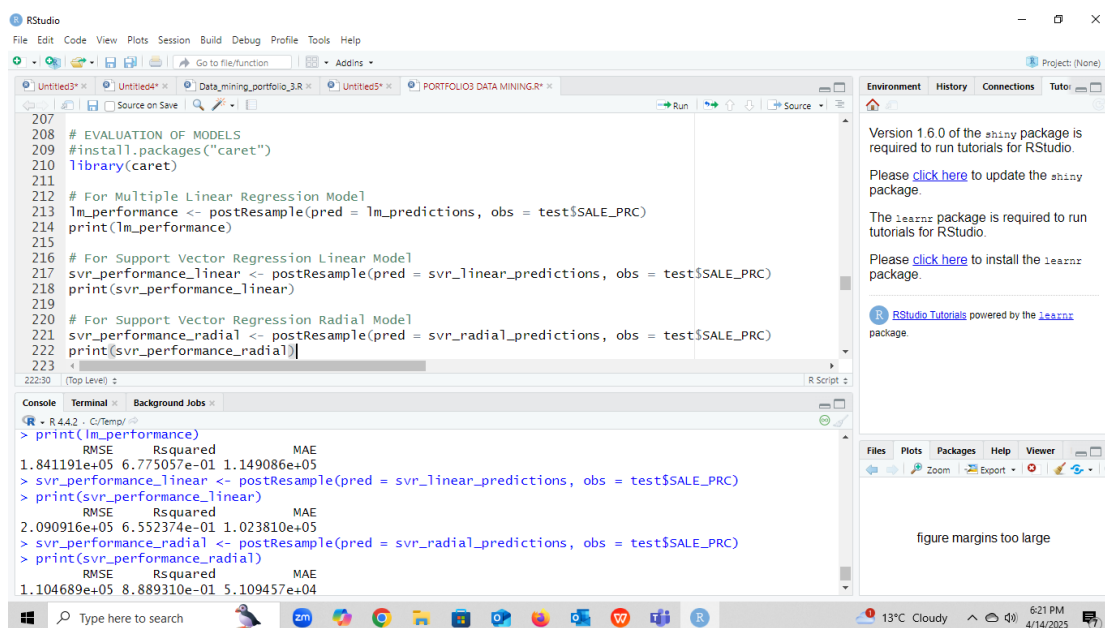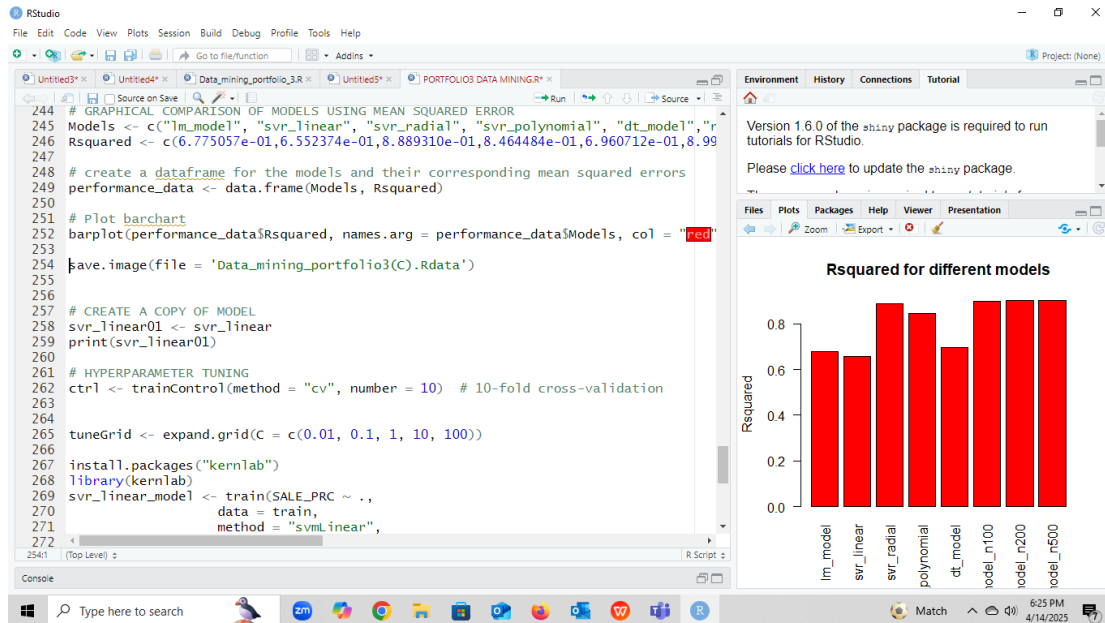


Fig. 19: Evaluation of models
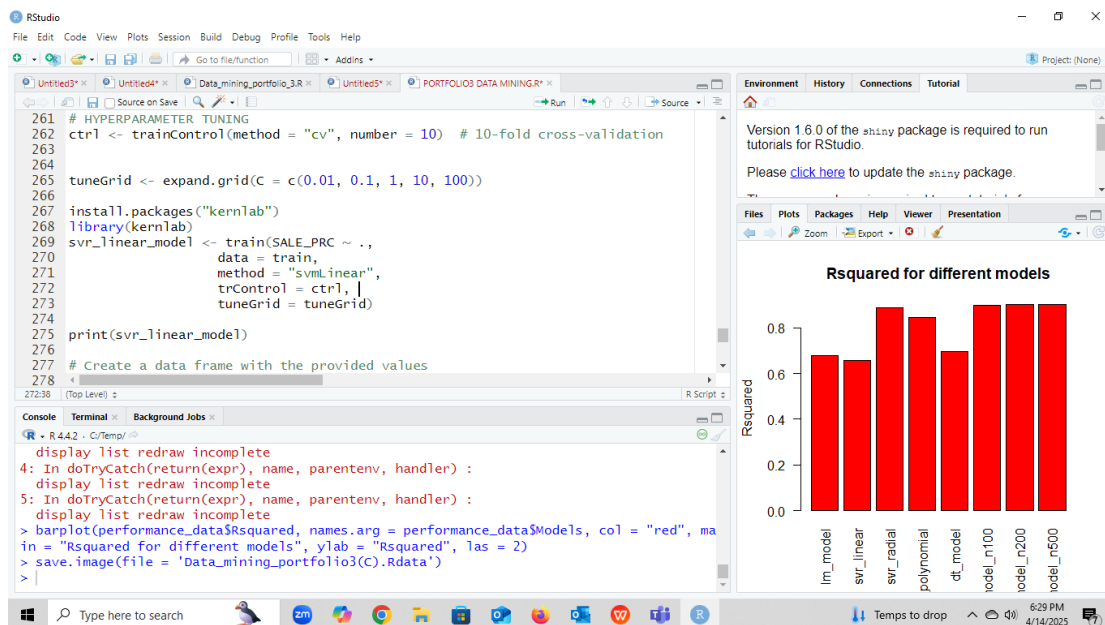
Fig. 20: Graphical Evaluation of models
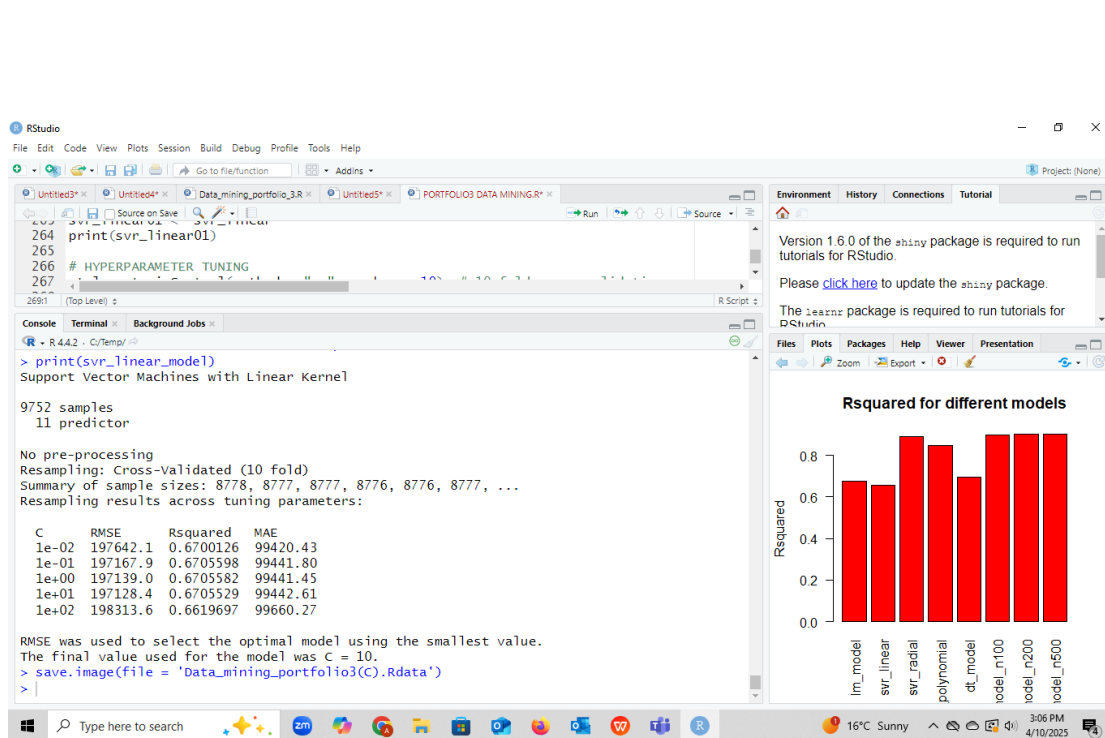


Fig. 21: Hyperparameter tuning

Fig. 22: Prediction with optimised support vector regression linear model

# CHAPTER 5

## DISCUSSION

The primary goal of this project is to predict house sale prices based on attributes such as size, distance from key areas, special features, age of property, and other relevant property characteristics. Predictive accuracy is key to supporting pricing decisions for buyers, sellers, and real estate professionals. The real estate industry thrives on accurate property valuation. In this project, eight (8) regression models were explored to predict house sale prices based on historical housing data. The models were evaluated using Mean Squared Error (MSE) as the primary metric, due to its sensitivity to large errors . Among the eight models developed, Support Vector Regression with a linear kernel emerged as the best performing model with the lowest MSE, making it the most suitable for business deployment. To ensure fair comparison, all models were train on the same training set and validated on the same test set for uniform and unbias evaluation. A vital phase in this project is the evaluation of all predictive models. The goal was not just to identify the most accurate model, but to select the one that would provide the most consistent and reliable predictions in a real-world house sale price setting. Hyperparameter was performed using grid search on support vector regression linear model to ensure the model performs at its best as shown in Fig…..

In the real estate industry, pricing is everything. A predictive model that accurately estimates house sale prices is not just a data science tool , it is a growth driver.

Deploying the support vector regression linear model would bring well defined business benefits such as risk reduction, revenue optimisation, efficiency in decision making, market intelligence and strategic planning, client confidence, promote stability across market changes among others.

# CHAPTER 6

## CONCLUSION

In conclusion, this project successfully demonstrates the application of regression modeling techniques in predicting house prices using machine learning algorithms. By following the structured Cross Industry Standard Process for Data Mining (CRISP-DM) methodology, each phase of the project from understanding the business problem to evaluating the predictive models was well addressed. A total of eight (8) machine learning models were developed and evaluated based on their performance in predicting house sale prices using a set of relevant explanatory variables. Among these, the support vector regression model proved to be the most effective, delivering superior accuracy compared to the other models with respect to having the least mean squared error. This outcome brings to light the potential of advanced regression techniques in real estate price prediction and validate the importance of careful model selection and evaluation in building reliable predictive systems.

# REFERENCES

Awad, M., Khanna, R., Awad, M. and Khanna, R., 2015. Support vector regression. *Efficient learning machines: Theories, concepts, and applications for engineers and system designers*, pp.67-80.

Basysyar, F.M. and Dwilestari, G., 2022. House price prediction using exploratory data analysis and machine learning with feature selection. *Acadlore Transactions on AI and Machine Learning*, *1*(1), pp.11-21.

Fahrmeir, L., Kneib, T., Lang, S. and Marx, B.D., 2022. Regression models. In *Regression: Models, methods and applications* (pp. 23-84). Berlin, Heidelberg: Springer Berlin Heidelberg

Foryś, I., 2022. Machine learning in house price analysis: regression models versus neural networks. *Procedia Computer Science*, *207*, pp.435-445.

Ho, W.K., Tang, B.S. and Wong, S.W., 2021. Predicting property prices with machine learning algorithms. *Journal of Property Research*, *38*(1), pp.48-70.

Kang, Y., Zhang, F., Peng, W., Gao, S., Rao, J., Duarte, F. and Ratti, C., 2021. Understanding house price appreciation using multi-source big geo-data and machine learning. *Land use policy*, *111*, p.104919.

Kuvalekar, A., Manchewar, S., Mahadik, S. and Jawale, S., 2020, April. House price forecasting using machine learning. In *Proceedings of the 3rd international conference on advances in science & technology (ICAST)*.

Manasa, J., Gupta, R. and Narahari, N.S., 2020, March. Machine learning based predicting house prices using regression techniques. In *2020 2nd International conference on innovative mechanisms for industry applications (ICIMIA)* (pp. 624-630). IEEE.

Monika, R., Nithyasree, J., Valarmathi, V., Hemalakshmi, G.R. and Prakash, N.B., 2021. House price forecasting using machine learning methods. *Turkish Journal of Computer and Mathematics Education*, *12*(11), pp.3624-3632.

Pardoe, I., 2020. *Applied regression modeling*. John Wiley & Sons.

Rigatti, S.J., 2017. Random forest. *Journal of Insurance Medicine*, *47*(1), pp.31-39.

Satish, G.N., Raghavendran, C.V., Rao, M.S. and Srinivasulu, C., 2019. House price prediction using machine learning. *Journal of Innovative Technology and Exploring Engineering*, *8*(9), pp.717-722.

Sharma, S., Arora, D., Shankar, G., Sharma, P. and Motwani, V., 2023, February. House price prediction using machine learning algorithm. In *2023 7th International Conference on Computing Methodologies and Communication (ICCMC)* (pp. 982-986). IEEE.

Truong, Q., Nguyen, M., Dang, H. and Mei, B., 2020. Housing price prediction via improved machine learning techniques. *Procedia Computer Science*, *174*, pp.433-442.

Yu, Y., Lu, J., Shen, D. and Chen, B., 2021. Research on real estate pricing methods based on data mining and machine learning. *Neural Computing and Applications*, *33*, pp.3925-3937.

Zhang, Z., 2016. Decision tree modeling using R. *Annals of translational medicine*, *4*(15), p.275.

Zulkifley, N.H., Rahman, S.A., Ubaidullah, N.H. and Ibrahim, I., 2020. House price prediction using a machine learning model: a survey of literature. *International Journal of Modern Education and Computer Science*, *12*(6), pp.46-54.