

# Machine Learning Master Project: California Housing Dataset

Citizen Data Science Program – Level 2

---

## Guidelines:

1. Answer each question with a clear and concise explanation of the Operator you used. Explain why you chose that particular operator and how it helps.
2. Outline the steps you followed to arrive at the solution.
3. Provide a screenshot of your process that demonstrated the relevant portion of your answer.
4. Name the document as the {project name - your name} before submitting it.
5. Ensure that you include both your answers word document and RapidMiner process as attachments in a reply to the person who sent you the assessment.
6. Write your name and email at the bottom of this page

## Project Description

This project consists of 16 questions related to your process in RapidMiner of the California Housing dataset. The questions are designed to assess your understanding of the key concepts and your ability to apply them in practical scenarios.

You will work on the California Housing dataset using RapidMiner to build a model of housing prices to predict median house values in California. The dataset contains information about housing prices in various districts of California. The dataset includes the following features:

1. longitude: Longitude value for the block in California, USA
2. latitude: Latitude value for the block in California, USA
3. housing\_median\_age: Median age of the house in the block
4. total\_rooms: Count of the total number of rooms (excluding bedrooms) in all houses in the block
5. total\_bedrooms: Count of the total number of bedrooms in all houses in the block
6. population: Count of the total number of population in the block
7. households: Count of the total number of households in the block
8. median\_income: Median of the total household income of all the houses in the block
9. ocean\_proximity: Type of the landscape of the block [ Unique Values : 'NEAR BAY', '<1H OCEAN', 'INLAND', 'NEAR OCEAN', 'ISLAND' ]
10. median\_house\_value: Median of the household prices of all the houses in the block

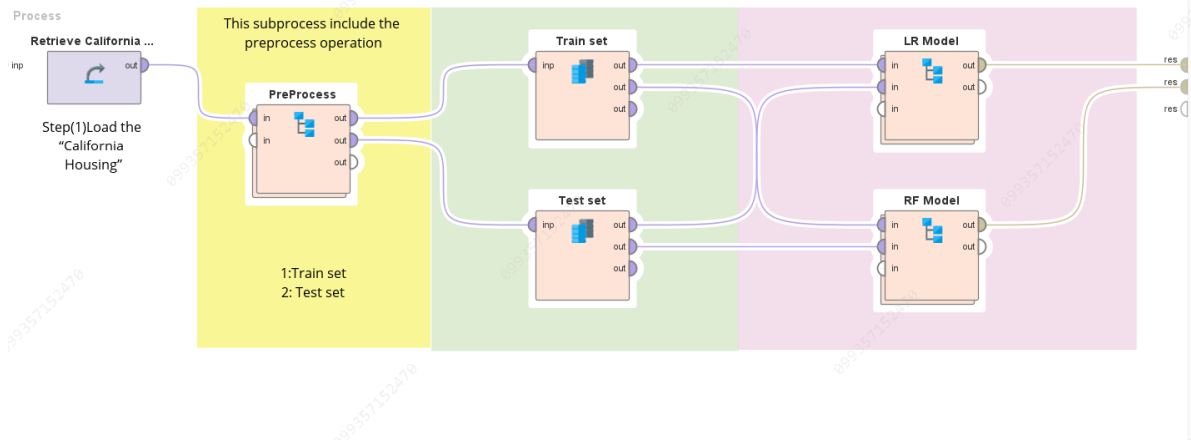
Name: Abdullah Mohammad Al Talaq

Email: [AlTalaqA@sabic.com](mailto:AlTalaqA@sabic.com)

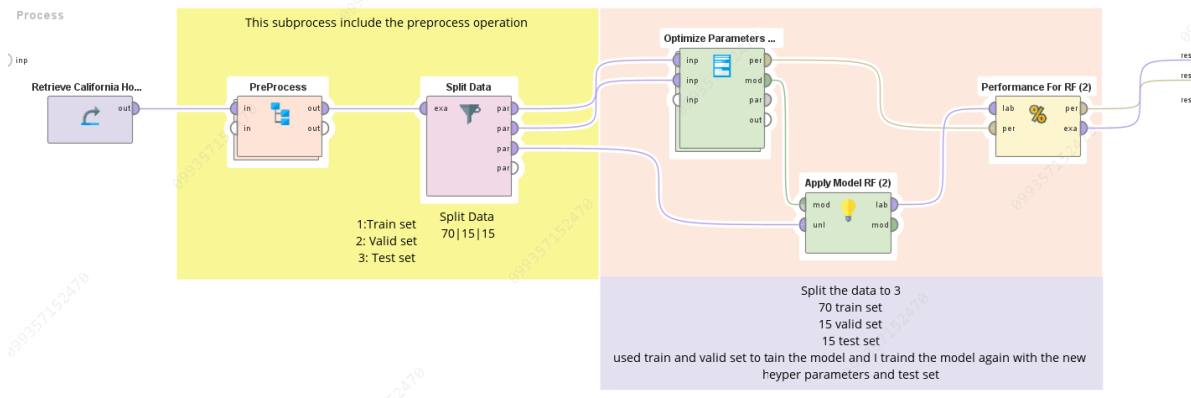
---

Section 0: The project includes three files.

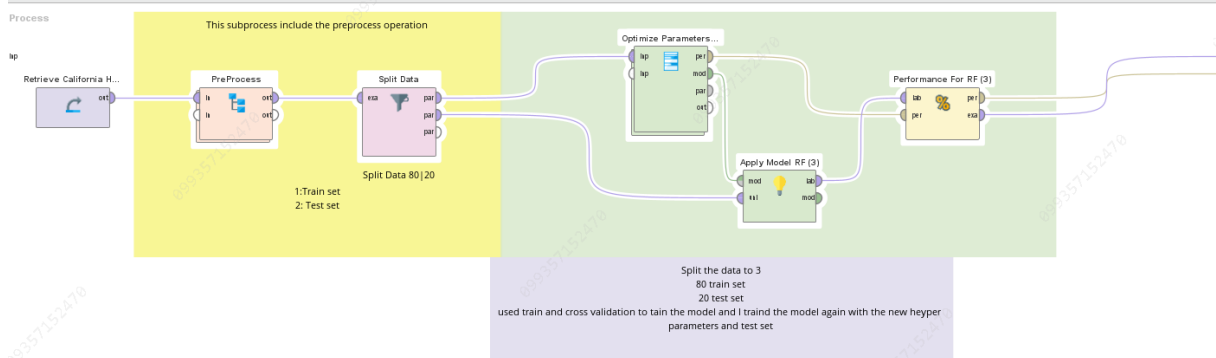
### 1- Machine Learning Process (Master)



### 2- Machine Learning Process (Master) Optimizations



### 3- Machine Learning Process (Master) Optimizations cross-validation



**Note there are sub process operations. To check you have to open the process in Ai studio (RapidMiner)**

Section 1: Data Exploration & Preprocessing

- 1. Preprocess the dataset by handling the missing values, if any, and. How many missing values are there? And how did you handle them?

total_bedrooms	Integer	207	Min 1	M... 6445	Average 537.871
----------------	---------	-----	----------	--------------	--------------------

Handling:

Views: DesignResultsTurbo PrepAuto ModelInteractive Analysis

Find data, operators...etcAll Studio

200%

Step(2) Replace null values

exa

exa

ori

pre

Step(2) Replace null values

Missing values were 207 in total\_bedrooms

Parameters

Step(2) Replace null values (Replace Missing Values)

attribute filter typeall

☐ invert selection

☐ include special attributes

defaultaverage

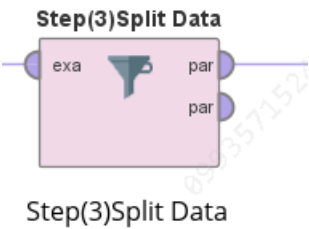
columnsEdit List (0)...

exa

Step(3) scaling

C

- 2. Split the dataset into training and testing sets using a 70:30 split ratio.



Edit Parameter List: partitions

Edit Parameter List: **partitions**  
The partitions that should be created.

ratio
0.7
0.3

3. Apply feature scaling to the dataset. Which scaling method did you use, and why this step is important?

100%

**Normalize**

exa exa  
ori  
pre

res  
res

feature scaling to the dataset

**Parameters**

**Normalize**

attribute filter type subset

attributes Select Attributes...

☐ invert selection

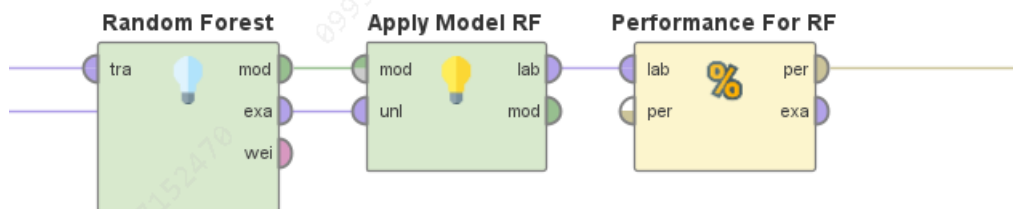
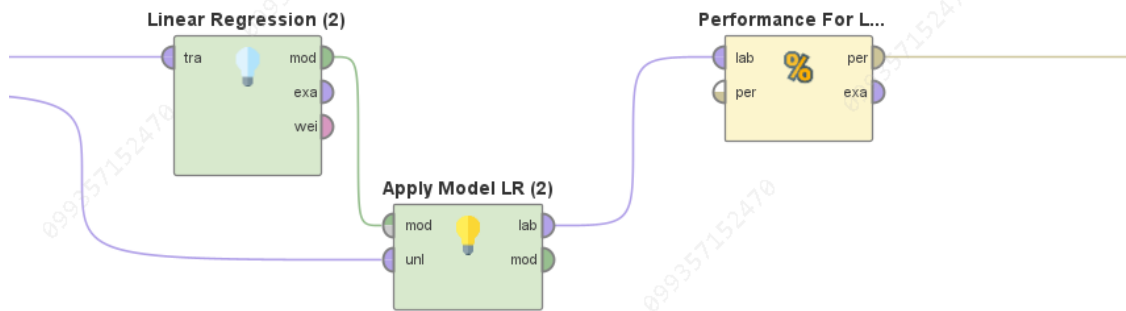
☐ include special attributes

method Z-transformation

Important process of standardizing data formats to ensure consistency and accuracy. By eliminating redundancies and errors, data normalization makes managing large volumes of information smoother and more efficient.

## Section 2: Model Selection &amp; Evaluation

1. Choose two algorithms from the given options: [Linear Regression, Random Forest, or Gradient Boosting Trees]. Justify your choice for each algorithm based on characteristics and suitability for the given dataset.



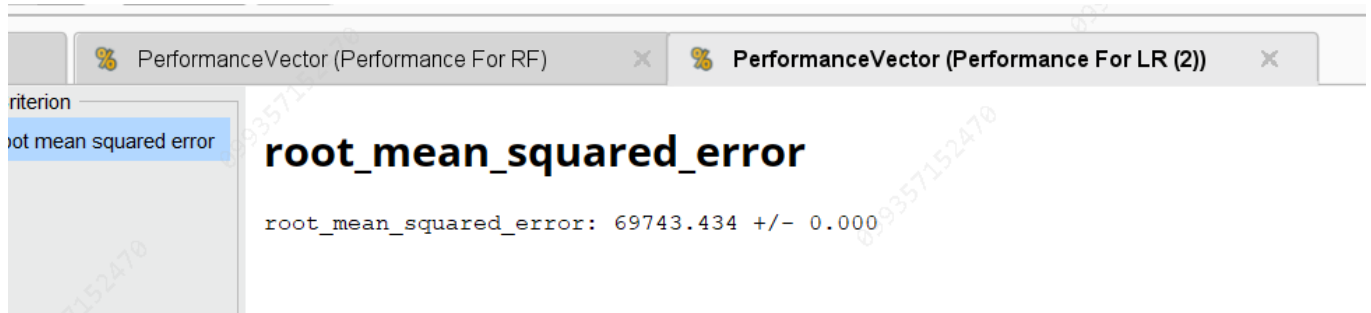
Justify for Linear: outcome variable is a numeric variable, and linear regression has fast and simplicity.

Justify for Random Forest: random forests models are not as sensitive to hyperparameters as other models, so it is easy to get a decent model up and running without having to do too much parameter tuning.

2. Train the first chosen algorithm on the training data to predict housing prices. Evaluate its performance using RMSE on testing data and explain what the RMSE performance metric is?

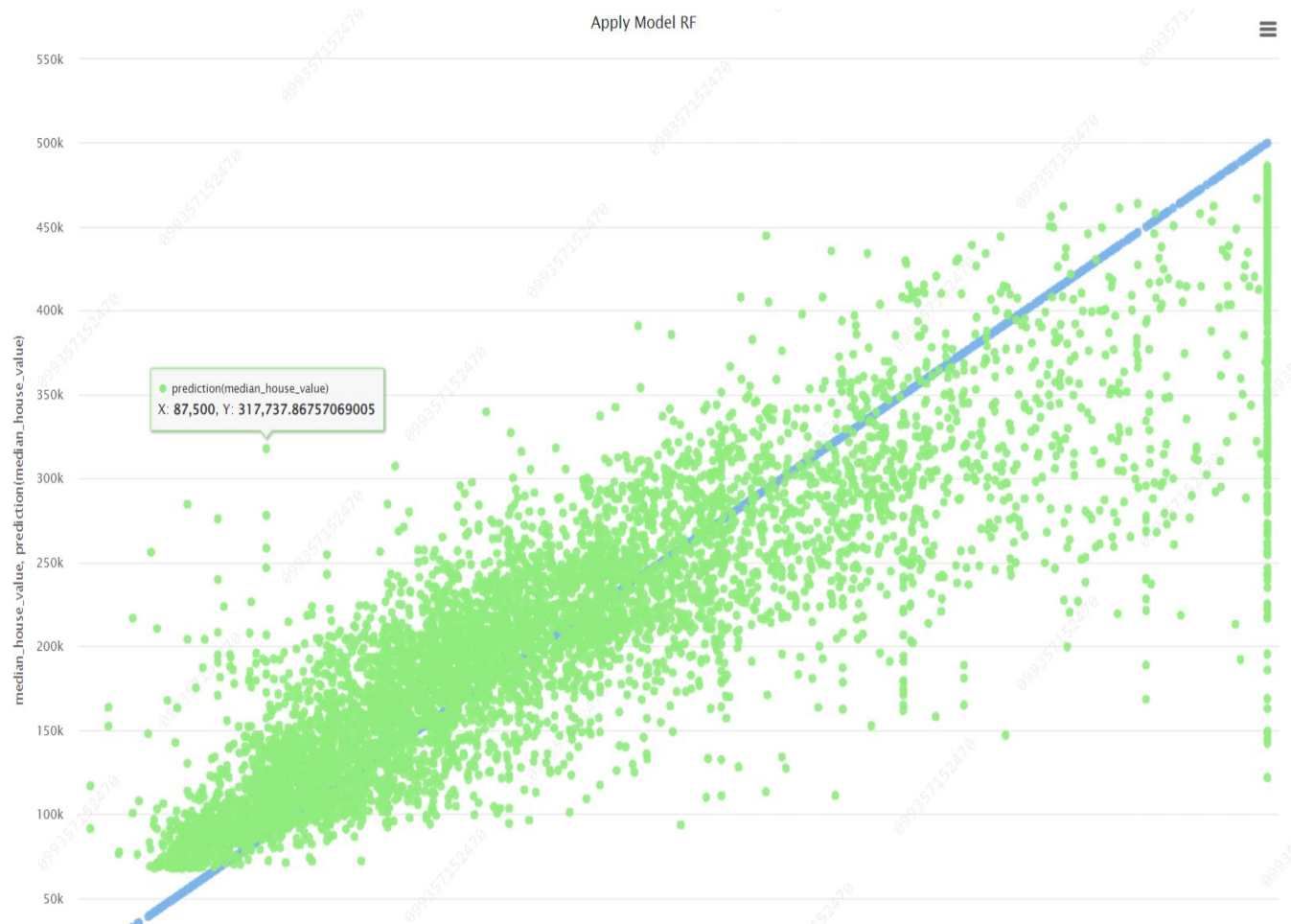


3. Train the second chosen algorithm on the training data and evaluate its performance on testing data (RMSE).



4. Compare the performance of the two chosen algorithms. Which algorithm performs better for predicting housing prices, and why?





Random forest performs better than Linear regression. RF predictions are close to actual values and you barely find outliers while in LR there are some predictions far from the range of y-axis for actual data. Also, there are prediction values less than zero, which is unlogical.

### Section 3: Hyperparameter Tuning & Optimization

1. Apply grid search or random search optimization to tune the hyperparameters of the better performing algorithm. Specify the hyperparameters you selected for optimization and explain their significance in the context of the chosen algorithm.

The screenshot displays the Orange3 data mining software interface. At the top, a workflow is visible with three main components: 'File Encoding', 'Set Role', and 'Optimize Parameters (Grid)'. The 'File Encoding' component has inputs for 'exa', 'ori', and 'pre'. The 'Set Role' component has inputs for 'exa' and 'ori', and is labeled 'Set House Price'. The 'Optimize Parameters (Grid)' component has inputs for 'inp' and 'inp', and outputs for 'per', 'mod', 'par', 'out', 'out', 'out', and 'out'. Below the workflow, there is a section for 'error handling' with a dropdown set to 'fail on error'. Other options include 'log performance' (checked), 'log all criteria' (unchecked), 'synchronize' (unchecked), and 'enable parallel execution' (checked).

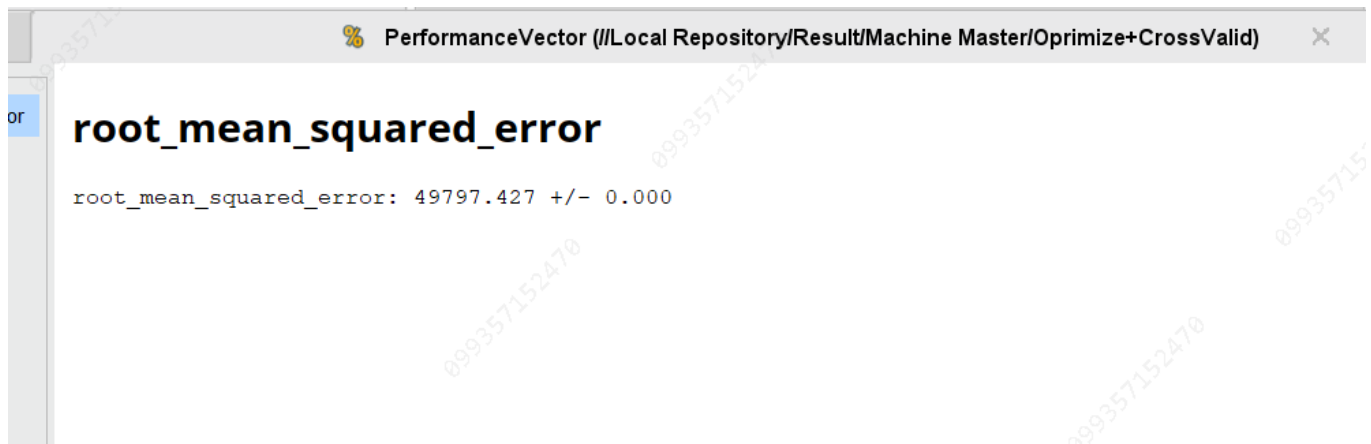
A dialog box titled 'Select Parameters: configure operator' is open in the foreground. It contains a list of operators on the left, including 'Step(3)Split Data (Split Data)', 'Test Set (Multiply)', 'Train Set (Multiply)', 'Random Forest (Random Forest)', 'Linear Regression (Linear Regression)', 'Apply Model LR (Apply Model)', and 'Performance For LR (Performance (Reg))'. The 'Random Forest (Random Forest)' operator is selected. The 'Parameters' section is empty. The 'Selected Parameters' section lists 'Random Forest.number\_of\_trees' and 'Random Forest.maximal\_depth'. Below this, there is a 'Grid/Range' section with fields for 'Min' (0.0), 'Max' (0.0), 'Steps' (0), and a 'Scale' dropdown set to 'linear'. At the bottom, there is a 'Value List' section with a large empty text area. The dialog also shows '2 parameters / 121 combinations selected' and buttons for 'OK' and 'Cancel'.

Number of trees: The main advantage of using number of trees that predictive performance tends to increase as the number of trees increases.

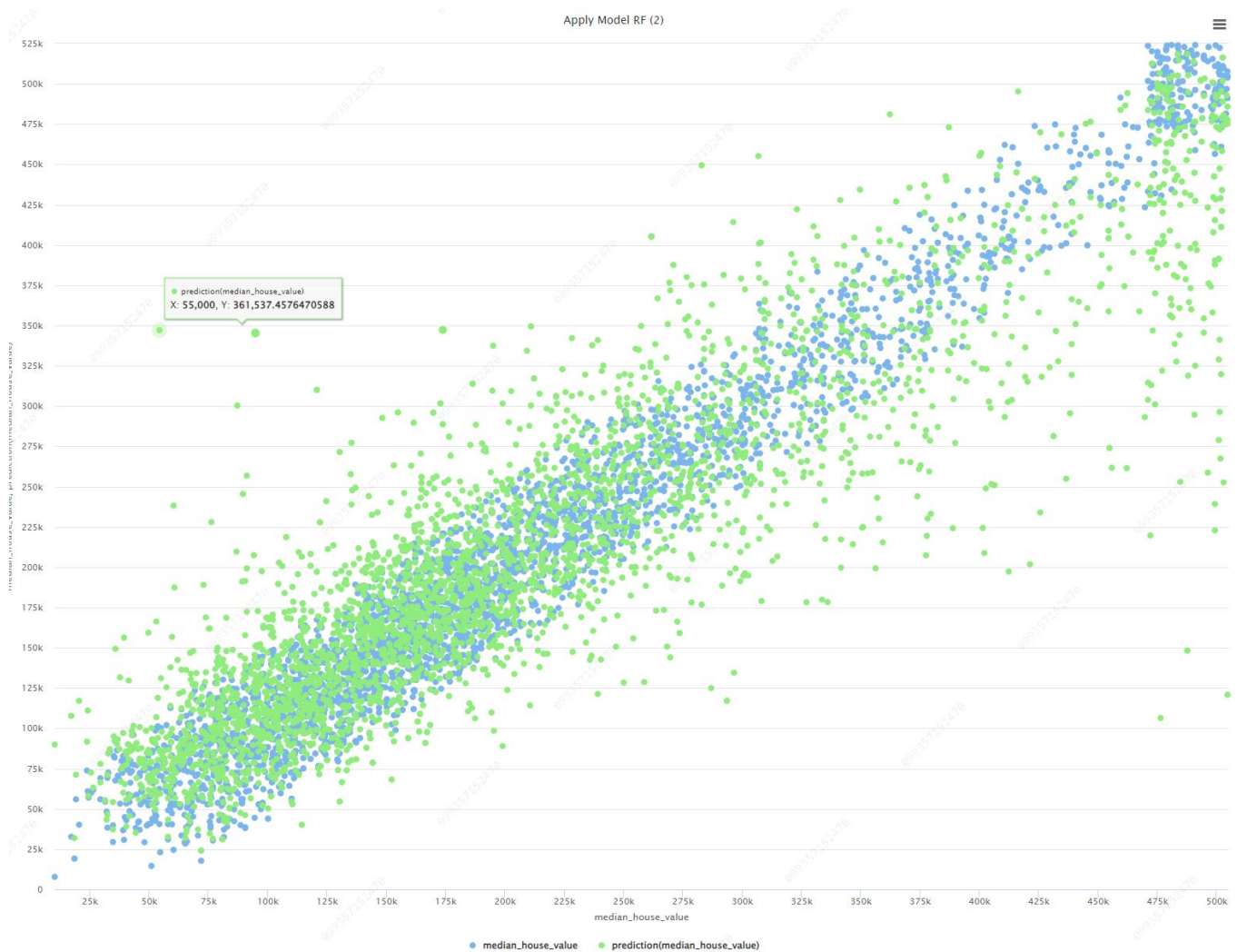
Maximal depth: Advantage of creating trees with a high max depth is that its predictive performance may benefit from doing so. In general, adding more splits to your trees will result in better classification.

2. Train the optimized model on the training data and evaluate its performance on testing data (RMSE).





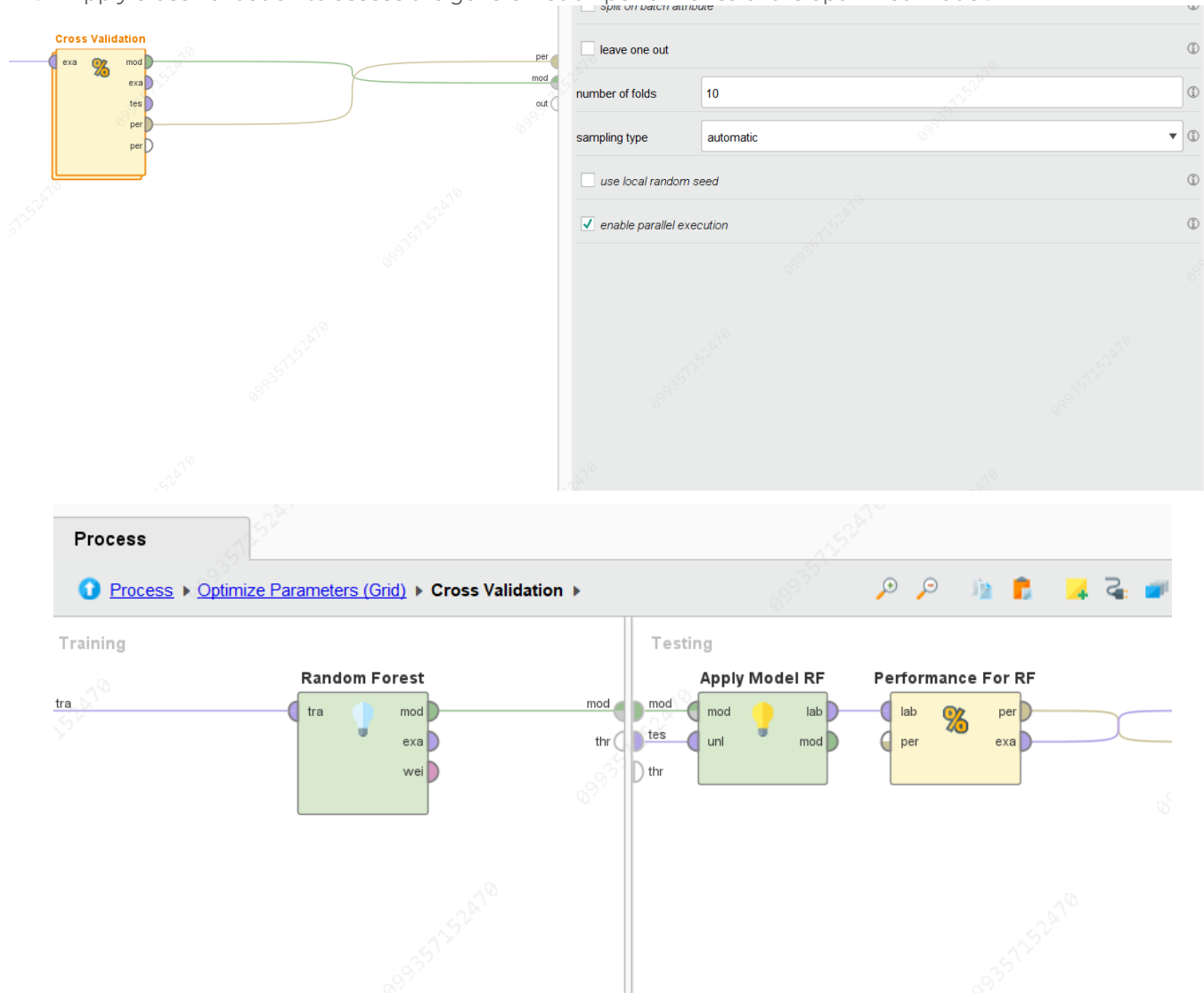
3. Show the performance of the optimized model with the previous model. Has the performance improved after optimization?



The performance becomes better, RMSE decreased.

## Section 4: Model Interpretation & Analysis

1. Apply cross-validation to assess the generalization performance of the optimized model.



## Section 5: Overfitting & Underfitting Analysis

1. Explain overfitting and underfitting in machine learning.

**overfitting:** A model gets trained with so much data, it starts learning from the noise and inaccurate data entries in our data set. And when testing with test data results in High variance. Then the model does not categorize the data correctly, because of too many details and noise.

**Underfitting:** A model is too simple to capture data complexities. It represents the inability of the model to learn the training data effectively result in poor performance both on the training and testing data.

## Section 6: Algorithm Strengths, Recommendation, and Implications

1. Summarize your findings and provide recommendation for predicting housing prices based on the provided dataset.

There should be an attribute for the city which the housing is in, having a data on the entirety of California is far more generalized to be modeled. According to google. California population is 39 million, and 485 cities. To make suitable model for such big city like this we need bigger dataset than this.

*\*Feel free to explore and experiment with the dataset and RapidMiner's capabilities beyond the questions mentioned above.\**

