

We split our data into two sets:

1. Trainset (X_train, y_train): This is the part of the data that we will use to train our model.
2. Testset (X_test, y_train): This part of the date is used to measure the quality of our model.

We use the function `train_test_split` to do that, and choose to put 30% of the data into the test set.

```
In [34]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=0)
```

Training the model

Here we choose to train a `DecisionTreeRegressor` model from library `sklearn`.

We choose the parameters of the model via cross validation on the training set.

The data used here is only the trainset.

```
In [35]: regressor = DecisionTreeRegressor(random_state=0)
parameters = {'max_depth': (3, 5, 7, 9, 10, 15, 20),
              'min_samples_split': (2, 4, 10, 20, 30, 50)}
model = GridSearchCV(regressor, parameters)
model.fit(X_train, y_train)
```

```
Out[35]:
GridSearchCV
  * estimator: DecisionTreeRegressor
    * DecisionTreeRegressor
```

Testing the model

Now the model is trained, we use the `predict` function to predict house prices using the variables.

We evaluate the quality of the model on the test set and also on the training set for comparison.

```
In [36]: predictTrain = model.predict(X_train)
predictTest = model.predict(X_test)
print("MSE train: (mean_squared_error(y_train, predictTrain)) | MSE test: (mean_squared_error(y_test, predictTest))")
print("MAE train: (mean_absolute_error(y_train, predictTrain)) | MAE test: (mean_absolute_error(y_test, predictTest))")
print("R squared train: (r2_score(y_train, predictTrain)) | R squared test: (r2_score(y_test, predictTest))")

MSE train: 27566745573.11343 | MSE test: 44642646644.47781
MAE train: 101759.05529247221 | MAE test: 123407.49951140535
R squared train: 0.7939614063620175 | R squared test: 0.6748342046837703
```

Our selected model performs relatively well being able to explain 67.5% of the variance in the data on the testing sample. However, the better performance of 79.4% variance explanation on the training set signals there is still some overfitting which may be addressed by using more regularization in the model or by gathering more data.

The two following figures represent the predicted values against the actual values and provide a visual way to assess the model quality. With a perfect model, every blue point would be on the red line. The further away from the red line they are, the bigger the error is.



Looking at the above plots we can visualize how the predictions of our model differ from the actual values.

Conclusion

We have performed some exploratory analysis on the house sales dataset and built a regression model to predict the price of the houses. Our proposed model could still be improved as the results show that it suffers from overfitting. Further exploration could consist of using different, more complex models.