

House Price Prediction with Linear Regression and Random Forest

The aim of this project is to predict real-estate prices using the machine learning

Loading the data

In the previous phase - Provisioning, data collection, data cleaning and data preparation were performed on a data scraped from website. In this file, this processed data will be used to train the models.

Modelling

In this stage, I decided to use several models and eventually I can decide which one performed the best in order to use in the next phase - Deployment. I will explore Linear, Lasso and Ridge Regression, Random Forest. All four will show different results for the accuracy. I decided to use these four models so as to check more features for comparing and different aspects.

I will compare the models by calculating the MAE, MSE, RMSE and the accuracy.

Evaluation

After training the models, it is important for the next steps to find out how good the performance of the models is. Based on this, it will be possible to conclude whether Modelling & Evaluation is successful or not.

Imports

```
import pandas as pd
import requests
import matplotlib.pyplot as plt
import seaborn as sns
from google.colab import files
from datetime import datetime
import io
import mpl_toolkits
import numpy as np
%matplotlib inline
```

```
# Load the data
local_file = files.upload()
train_data = io.BytesIO(local_file['data.csv'])
df = pd.read_csv(train_data)
```



Choose Files

No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving data.csv to data.csv

➤ Preparing the data for training the models

Train-Test Split dataset

Necessary imports

```
from sklearn.linear_model import LinearRegression, Ridge, Lasso
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import cross_val_score, train_test_split, GridSearchCV
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 96 entries, 0 to 95
Data columns (total 10 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Unnamed: 0      96 non-null    int64
1   AGENCY          96 non-null    float64
2   TYPE            96 non-null    float64
3   STREET NAME     96 non-null    object
4   POSTCODE        96 non-null    int64
5   LIVING_AREA     96 non-null    float64
6   ROOMS           96 non-null    float64
7   LONGITUDE       96 non-null    float64
8   LATITUDE        96 non-null    float64
9   PRICE_LOG       96 non-null    float64
dtypes: float64(7), int64(2), object(1)
memory usage: 7.6+ KB
```

```
df.drop('Unnamed: 0', axis = 1, inplace = True)
```

```
df.isnull().sum()
```

```

AGENCY      0
TYPE        0
STREET_NAME 0
POSTCODE    0
LIVING_AREA 0
ROOMS       0
LONGITUDE   0
LATITUDE    0
PRICE_LOG   0
dtype: int64

```

Analyzing the numeric features.

```
numeric_features = df.select_dtypes(include=[np.number])
```

```
numeric_features.columns
```

```

Index(['AGENCY', 'TYPE', 'POSTCODE', 'LIVING_AREA', 'ROOMS', 'LONGITUDE',
      'LATITUDE', 'PRICE_LOG'],
      dtype='object')

```

```
# set the target and predictors
```

```
y = df.PRICE_LOG # target
```

```
# use only those input features with numeric data type
```

```
df_temp = df.select_dtypes(include=["int64","float64"])
```

```
X = df_temp.drop(["PRICE_LOG"],axis=1) # predictors
```

To split the dataset, I will use random sampling with 80/20 train-test split; that is, 80% of the dataset will be used for training and set aside 20% for testing:

X_train and y_train contain data for the training model. X_test and y_test contain data for the testing model. X and y are features and target variable names.

```
# split the dataset into train and test sets
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=0)
```

▼ Modelling

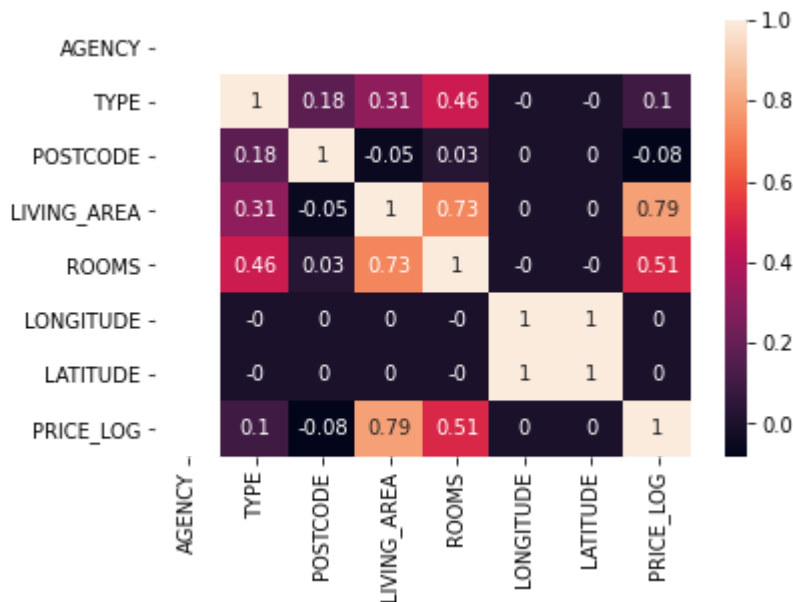
Two models will be built and evaluated by their performances with R-squared metric. Additionally, insights on the features that are strong predictors of house prices, will be analysed .

Linear Regression

To fit a linear regression model, the features which have a high correlation with the target variable PRICE are selected. By looking at the correlation matrix, it is noticeable that the rooms and the living area have a strong correlation with the price ('Log of price').

```
correlation_matrix = df.corr().round(2)
# annot = True to print the values inside the square
sns.heatmap(data=correlation_matrix, annot=True)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f8c650a0d90>



Fitting models describes the relationship between a response variable and one or more predictor variables.

```
lr = LinearRegression()
# fit optimal linear regression line on training data
lr.fit(X_train,y_train)
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)

LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)

coeff_df = pd.DataFrame(lr.coef_,X.columns,columns=['Coefficient'])
coeff_df
```

	Coefficient
AGENCY	0.000000e+00
TYPE	-9.364483e-02
POSTCODE	2.009291e-03
LIVING AREA	1.015017e-02

What coefficient of data says:

Holding all other features fixed, a 1 unit increase in Type is associated with a decrease of 0.20 euros. Holding all other features fixed, a 1 unit increase in Postcode is associated with a decrease of 0.20 euros. Holding all other features fixed, a 1 unit increase in Living area is associated with an increase of 0.01 euros . Holding all other features fixed, a 1 unit increase in Rooms is associated with an increase of 0.01 euros .

Root Mean Square Error (RMSE) is the standard deviation of the residuals (prediction errors). Residuals are a measure of how far from the regression line data points are; RMSE is a measure of how spread out these residuals are. In other words, it tells you how concentrated the data is around the line of best fit.

RMSE values between 0.2 and 0.5 shows that the model can relatively predict the data accurately.

```
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error
```

```
#predict y_values using X_test set
yr_hat = lr.predict(X_test)
```

Comparing these metrics:

MAE is the easiest to understand because it's the average error. MSE is more popular than MAE because MSE "punishes" larger errors, which tends to be useful in the real world. RMSE is even more popular than MSE because RMSE is interpretable in the "y" units.

```
# model evaluation for testing set
print('MAE:', mean_absolute_error(y_test, yr_hat))
print('MSE:', mean_squared_error(y_test, yr_hat))
print('RMSE:', np.sqrt(mean_squared_error(y_test, yr_hat)))
```

MAE: 0.3410280435512464

MSE: 0.13885727231548534
 RMSE: 0.37263557575127654

In this case, the RMSE vary between 0.2 and 0.3 - changing depending on the website information, so it is relatively accurate.

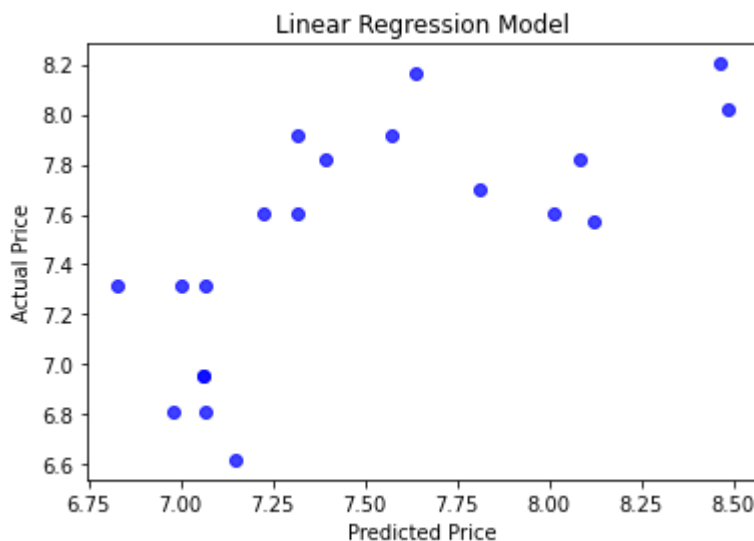
Accuracy is one metric for evaluating classification models. Accuracy is the fraction of predictions our model got right.

```
lr_score =lr.score((X_test),y_test)
print("Accuracy: ", lr_score)
```

Accuracy: 0.3480141020654395

The model showed accuracy of 83% after its training.

```
actual_values = y_test
plt.scatter(yr_hat, actual_values, alpha=.75,
            color='b') #alpha helps to show overlapping data
plt.xlabel('Predicted Price')
plt.ylabel('Actual Price')
plt.title('Linear Regression Model')
plt.show()
#pltrandom_state=None.show()
```



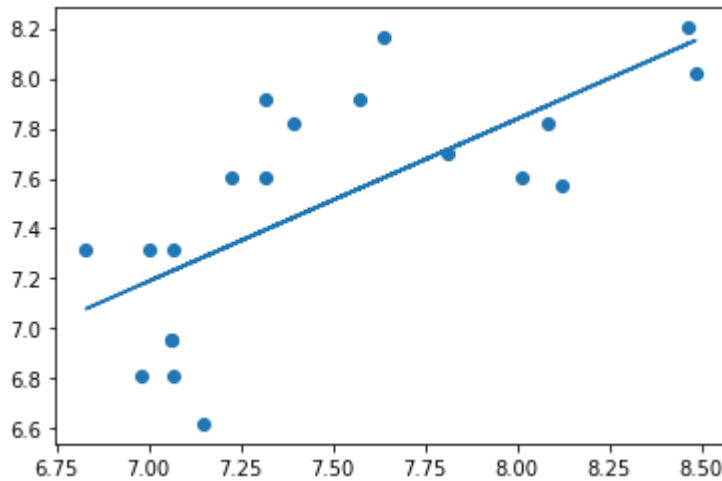
```
from scipy import stats
```

```
#Execute a method that returns the important key values of Linear Regression
slope, intercept, r, p, std_err = stats.linregress(yr_hat, y_test)
```

```
#Create a function that uses the slope and intercept values to return a new value. This new v
```

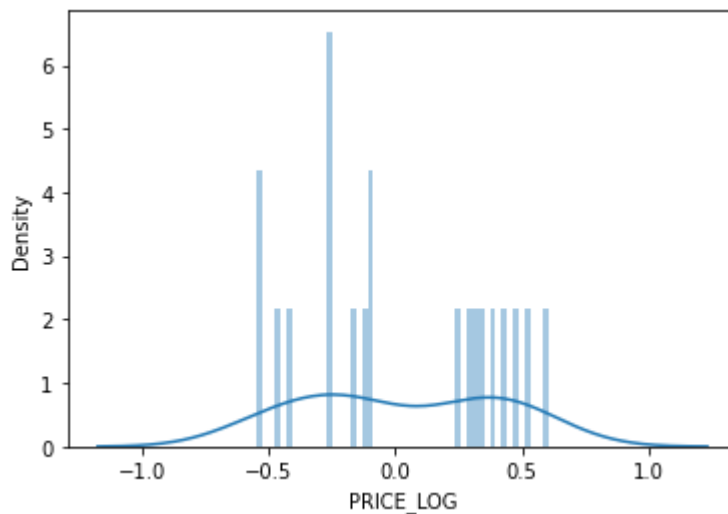
```
#Create a function that uses the slope and intercept values to return a new value. This new v
def myfunc(x):
    return slope * x + intercept

mymodel = list(map(myfunc, yr_hat))
#Draw the scatter plot
plt.scatter(yr_hat, y_test)
#Draw the line of linear regression
plt.plot(yr_hat, mymodel)
plt.show()
```



```
sns.distplot((y_test-yr_hat),bins=50);
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `di
warnings.warn(msg, FutureWarning)
```



In the above histogram plot, the data is in bell shape (Normally Distributed), which means the model has done good predictions.

Using cross-validation to see whether the model is over-fitting the data.

```
# cross validation to find 'validate' score across multiple samples, automatically does Kfold
lr_cv = cross_val_score(lr, X, y, cv = 5, scoring= 'r2')
print("Cross-validation results: ", lr_cv)
print("R2: ", lr_cv.mean())
```

```
Cross-validation results: [0.34270636 0.62650251 0.16944591 0.49666958 0.41319454]
R2: 0.40970377784855394
```

It doesn't appear that for this train-test dataset the model is over-fitting the data (the cross-validation performance is very close in value).

Regularization:

The alpha parameter in ridge and lasso regularizes the regression model. The regression algorithms with regularization differ from linear regression in that they try to penalize those features that are not significant in our prediction. Ridge will try to reduce their effects (i.e., shrink their coefficients) in order to optimize all the input features. Lasso will try to remove the not-significant features by making their coefficients zero. In short, Lasso (L1 regularization) can eliminate the not-significant features, thus performing feature selection while Ridge (L2 regularization) cannot.

Lasso regression

```
lasso = Lasso(alpha = 1) # sets alpha to almost zero as baseline
lasso.fit(X_train, y_train)
```

```
Lasso(alpha=1, copy_X=True, fit_intercept=True, max_iter=1000, normalize=False,
      positive=False, precompute=False, random_state=None, selection='cyclic',
      tol=0.0001, warm_start=False)
```

RMSE tells you how concentrated the data is around the line of best fit.

```
# model evaluation for training set
y_train_l_predict = lasso.predict(X_train)
rmse = (np.sqrt(mean_squared_error(y_train, y_train_l_predict)))
```

```
print("The model performance for training set:")
print('RMSE is {}'.format(rmse))
```

```
The model performance for training set:
RMSE is 0.32376161625960365
```


RMSE values between 0.2 and 0.5 shows that the model can relatively predict the data accurately. In this case, it is 0.5, so it is relatively accurate.

```
# model evaluation for testing set
y_test_l_predict = lasso.predict(X_test)
rmse = (np.sqrt(mean_squared_error(y_test, y_test_l_predict)))
print("The model performance for testing set:")
print('RMSE is {}'.format(rmse))
```

```
The model performance for testing set:
RMSE is 0.3609516473883934
```

RMSE values between 0.2 and 0.5 shows that the model can relatively predict the data accurately. In this case, it is 0.5, so it is relatively accurate.

```
#predict y_values using X_test set
yr_lasso = lasso.predict(X_test)
```

Accuracy is one metric for evaluating classification models. Accuracy is the fraction of predictions our model got right.

```
lasso_score =lasso.score((X_test),y_test)
print("Accuracy: ", lasso_score)
```

```
Accuracy:  0.388258943135693
```

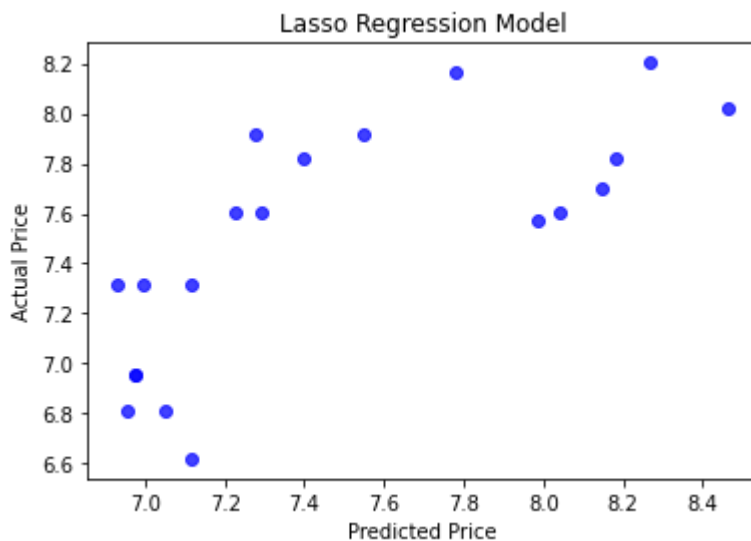
The model showed accuracy of 65% after its training.

```
lasso_cv = cross_val_score(lasso, X, y, cv = 5, scoring = 'r2')
print ("Cross-validation results: ", lasso_cv)
print ("R2: ", lasso_cv.mean())
```

```
Cross-validation results:  [0.43906093 0.69386751 0.22888875 0.56034884 0.37990456]
R2:  0.46041412107895907
```

```
actual_values = y_test
plt.scatter(yr_lasso, actual_values, alpha=.75,
            color='b') #alpha helps to show overlapping data
plt.xlabel('Predicted Price')
plt.ylabel('Actual Price')
plt.title('Lasso Regression Model')
plt.show()
```

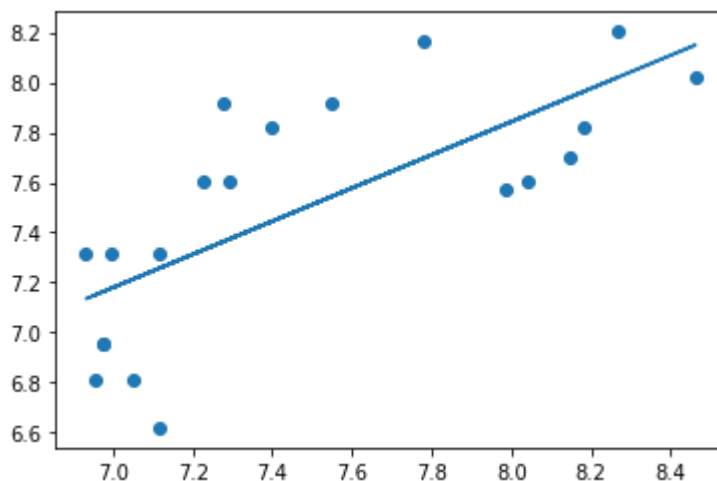
```
#plt.random_state=None.show()
```



```
from scipy import stats
```

```
#Execute a method that returns the important key values of Linear Regression
slope, intercept, r, p, std_err = stats.linregress(yr_lasso, y_test)
#Create a function that uses the slope and intercept values to return a new value. This new v
def myfunc(x):
    return slope * x + intercept
```

```
mymodel = list(map(myfunc, yr_lasso))
#Draw the scatter plot
plt.scatter(yr_lasso, y_test)
#Draw the line of linear regression
plt.plot(yr_lasso, mymodel)
plt.show()
```



Ridge regression

```
ridge = Ridge(alpha = 1) # sets alpha to a default value as baseline
```

```

ridge.fit(X_train, y_train)

Ridge(alpha=1, copy_X=True, fit_intercept=True, max_iter=None, normalize=False,
       random_state=None, solver='auto', tol=0.001)

# model evaluation for training set
y_train_r_predict = ridge.predict(X_train)
rmse = (np.sqrt(mean_squared_error(y_train, y_train_r_predict)))

print("The model performance for training set:")
print('RMSE is {}'.format(rmse))

The model performance for training set:
RMSE is 0.3099897204081143

```

RMSE values between 0.2 and 0.5 shows that the model can relatively predict the data accurately. In this case, it is 0.2, so it is relatively accurate.

```

# model evaluation for testing set
y_test_r_predict = ridge.predict(X_test)
rmse = (np.sqrt(mean_squared_error(y_test, y_test_r_predict)))
print("The model performance for testing set:")
print('RMSE is {}'.format(rmse))

The model performance for testing set:
RMSE is 0.3727507384589643

```

RMSE values between 0.2 and 0.5 shows that the model can relatively predict the data accurately. In this case, it is rounded to 0.2, so it is relatively accurate.

```

#predict y_values using X_test set
yr_ridge = ridge.predict(X_test)

```

Accuracy is one metric for evaluating classification models. Accuracy is the fraction of predictions our model got right.

```

ridge_score =ridge.score((X_test),y_test)
print("Accuracy: ", ridge_score)

```

```
Accuracy: 0.34761104841913404
```

The model showed accuracy of 83% after its training.

```

ridge_cv = cross_val_score(ridge, X, y, cv = 5, scoring = 'r2')
print ("Cross-validation results: ", ridge_cv)
print ("R2: ", ridge_cv.mean())

```

```

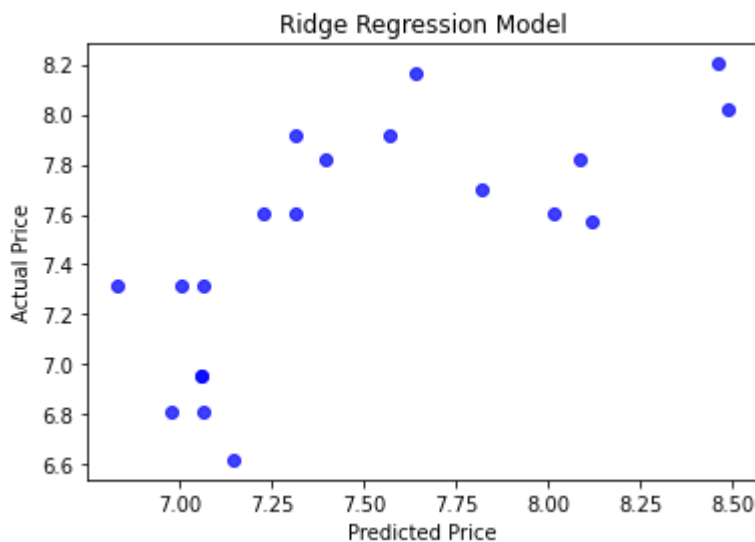
Cross-validation results: [0.3440135  0.62904074 0.16941042 0.49693456 0.41267706]
R2:  0.41041525616156677

```

```

actual_values = y_test
plt.scatter(yr_ridge, actual_values, alpha=.75,
            color='b') #alpha helps to show overlapping data
plt.xlabel('Predicted Price')
plt.ylabel('Actual Price')
plt.title('Ridge Regression Model')
plt.show()
#pltrandom_state=None.show()

```



```

from scipy import stats

```

```

#Execute a method that returns the important key values of Linear Regression
slope, intercept, r, p, std_err = stats.linregress(yr_ridge, y_test)
#Create a function that uses the slope and intercept values to return a new value. This new v
def myfunc(x):
    return slope * x + intercept

mymodel = list(map(myfunc, yr_ridge))
#Draw the scatter plot
plt.scatter(yr_ridge, y_test)
#Draw the line of linear regression
plt.plot(yr_ridge, mymodel)
plt.show()

```



13/18

```
# import the model we are using
from sklearn.ensemble import RandomForestRegressor
# Instantiate model with 1000 decision trees
rf = RandomForestRegressor(n_estimators = 1000, random_state = 42)
# Train the model on training data
rf.fit(X_train, y_train)

RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
                        max_depth=None, max_features='auto', max_leaf_nodes=None,
                        max_samples=None, min_impurity_decrease=0.0,
                        min_impurity_split=None, min_samples_leaf=1,
                        min_samples_split=2, min_weight_fraction_leaf=0.0,
                        n_estimators=1000, n_jobs=None, oob_score=False,
                        random_state=42, verbose=0, warm_start=False)
```

Making predictions on the test set:

When performing regression, the absolute error should be used. It needs to be checked how far away the average prediction is from the actual value so the absolute value has to be calculated.

```
# Use the forest's predict method on the test data
predictions = rf.predict(X_test)
# Calculate the absolute errors
errors = abs(predictions - y_test)
# Print out the mean absolute error (mae)
print('Mean Absolute Error:', round(np.mean(errors), 2), 'degrees.')
```

Mean Absolute Error: 0.29 degrees.

There is a 0.16 improvement.

Determine performance metrics:

To put the predictions in perspective, accuracy can be calculated by using the mean average percentage error subtracted from 100 %.

Accuracy is one metric for evaluating classification models. Accuracy is the fraction of predictions our model got right.

```
# Calculate mean absolute percentage error (MAPE)
mape = 100 * (errors / y_test)
# Calculate and display accuracy
accuracy = 100 - np.mean(mape)
print('Accuracy:', round(accuracy, 2), '%.')
```

Accuracy: 96.08 %.

The model has learned how to predict the price with 98% accuracy.

```
rfr = RandomForestRegressor()
rfr.fit(X_train, y_train) # gets the parameters for the rfr model
rfr_cv = cross_val_score(rfr,X, y, cv = 5, scoring = 'r2')
print("R2: ", rfr_cv.mean())
```

```
R2: 0.4990869758300721
```

The performance of Random forest is slightly better than the Linear regression. The model parameters can be optimised for better performance using gridsearch.

```
#Random forest determined feature importances
rfr.feature_importances_

array([0.          , 0.02773612, 0.09214652, 0.84844164, 0.03167572,
        0.          , 0.          ])
```

▼ Plotting the Feature Importance

Finding the features that are the most promising predictors:

```
importance = rfr.feature_importances_

# map feature importance values to the features
feature_importances = zip(importance, X.columns)

#list(feature_importances)
sorted_feature_importances = sorted(feature_importances, reverse = True)

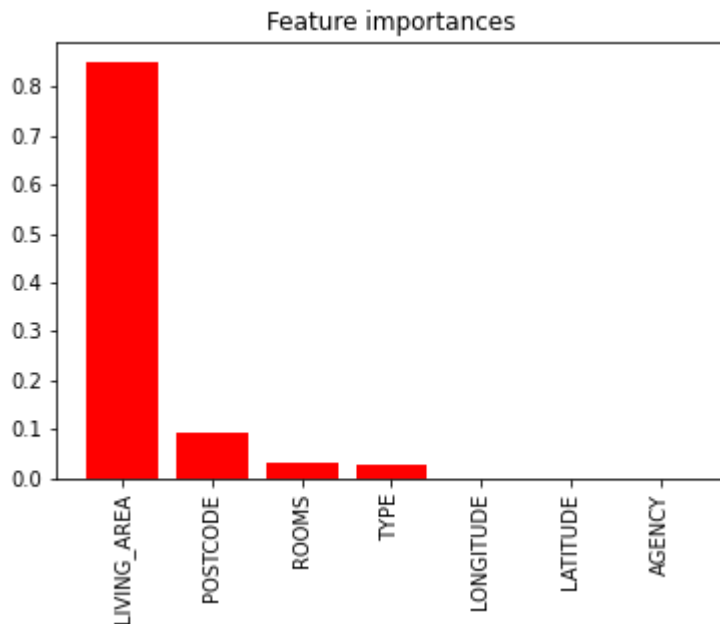
#print(sorted_feature_importances)
top_15_predictors = sorted_feature_importances[0:15]
values = [value for value, predictors in top_15_predictors]
predictors = [predictors for value, predictors in top_15_predictors]
print(predictors)

['LIVING_AREA', 'POSTCODE', 'ROOMS', 'TYPE', 'LONGITUDE', 'LATITUDE', 'AGENCY']
```

Plotting the feature importance of the Random forest.

Plotting the feature importances to illustrate the disparities in the relative significance of the

```
plt.figure()
plt.title( "Feature importances")
plt.bar(range(len(predictors)), values,color="r", align="center");
plt.xticks(range(len(predictors)), predictors, rotation=90);
```



The idea behind the plotting of feature importance is that after evaluating the performance of the model, the values of a feature of interest must be permuted and reevaluate model performance. The feature importance (variable importance) describes which features are relevant.

Random Forest determined that overall the living area of a home is by far the most important predictor. Following are the sizes of above rooms and postcode.

```
# tune models GridsearchCV
from sklearn.model_selection import GridSearchCV
parameters = {'n_estimators':range(10,300,10), 'criterion':('mse','mae'), 'max_features':('au

gs = GridSearchCV(rfr,parameters,scoring='neg_mean_absolute_error',cv=3)
gs.fit(X_train,y_train)

gs.best_score_
gs.best_estimator_

# test ensembles
tpred_lm = lr.predict(X_test)
tpred_lml = lasso.predict(X_test)
tpred_rf = gs.best_estimator_.predict(X_test)

from sklearn.metrics import mean_absolute_error
print(mean_absolute_error(y_test,tpred_lm))
```



```
print(mean_absolute_error(y_test, tpred_lm1))
print(mean_absolute_error(y_test, tpred_rf))

print(mean_absolute_error(y_test, (tpred_lm+tpred_rf)/2))

11.979225350971744
0.7719972599914945
0.25448917183278896
5.995337038314934
```

Saving the model with the best performance to be used in the deployment part of the project.

```
# Saving the model
import pickle
filename = 'classifier.pkl'
pickle.dump(rfr, open(filename, 'wb'))
```

Conclusion

I used four models to determine the accuracy - Linear Regression, Lasso Regression and Ridge Regression, Random Forest.

From the exploring of the models RMSE:

- Linear Regression score: 0.2003 (0.1887)
- Lasso score: 0.5 (0.4675)
- Ridge score: 0.2 (0.1877)
- Random forest score: 0.2372

RMSE values between 0.2 and 0.5 shows that the model can relatively predict the data accurately. All of the models showed values in this range.

From the exploring of the models accuracy:

- Linear Regression score: 0.80 (80%)
- Lasso score: 0.82 (82%)
- Ridge score: 0.86 (86%)
- Random forest score: 98.13 %

From the exploring of the models cross-validation:

- Linear Regression score: R2: 0.7308604883584712

- Lasso score: R2: 0.6532616143265344
- Ridge score: R2: 0.7310756447849953
- Random forest: R2: 0.7742740242196954

Random forest turns out to be the more accurate model for predicting the house price.

All of the models showed RMSE values between 0.2 and 0.5 so that they show relatively accurate predictions of the data.

I evaluated the models performances with R-squared metric and the one that is overfitting the least is the Linear Regression.

In the end, I tried three different models and evaluated them using Mean Absolute Error. I chose MAE because it is relatively easy to interpret and outliers aren't particularly bad in for this type of model. The one I will be using for the deployment is the **Random forest**.